# ALICE Multi-site Data Transfer Tests on a Wide Area Network

S. Bagnasco[a], R. Barbera[b,*], F. Carminati[c], P. Cerello[a], D. Di Bari[d], G. Donvito[e], E. Fragiacomo[f], A. Fritz[g], G. Lo Re[h], M. Luvisetto[i], M. Masera[j], F. Minafra[d], D. Mura[k], M. Sitta[a], J. Švec[l], R. Turrisi[m]

[a] *Istituto Nazionale di Fisica Nucleare, Sezione di Torino – ITALY*
[b] *Dipartimento di Fisica e Astronomia dell'Università di Catania and Istituto Nazionale di Fisica Nucleare, Sezione di Catania – ITALY*
[c] *CERN, Geneva – SWITZERLAND*
[d] *Dipartimento IA di Fisica dell?università di Bari and Istituto Nazionale di Fisica Nucleare, Sezione di Bari – ITALY*
[e] *Istituto Nazionale di Fisica Nucleare, Sezione di Bari – ITALY*
[f] *Istituto Nazionale di Fisica Nucleare, Sezione di Trieste – ITALY*
[g] *University of Houston – USA*
[h] *Istituto Nazionale di Fisica Nucleare, CNAF, Bologna  – ITALY*
[i] *Dipartimento di Fisica dell'Università di Bologna and Istituto Nazionale di Fisica Nucleare, Sezione di Bologna – ITALY*
[j] *Dipartimento di Fisica Sperimentale dell'Università di Torino and Istituto Nazionale di Fisica Nucleare, Sezione di Torino – ITALY*
[k] *Istituto Nazionale di Fisica Nucleare, Sezione di Cagliari – ITALY*
[l] *Institute of Physic, Academy of Sciences, Prague – CZECH REPUBLIC*
[m] *Istituto Nazionale di Fisica Nucleare, Sezione di Padova – ITALY*

## Abstract

*Next generation high energy physics experiments planned at the CERN Large Hadron Collider have very demanding needs in terms of computing power, mass storage, and network bandwidth. In this paper we report on a series of realistic multi-site data transfer tests on a wide area network performed within the ALICE Experiment in order to spot possible bottlenecks and pin down critical elements and parameters of actual research networks.*

## 1. Introduction

Next generation high energy physics experiments planned at the CERN [1] Large Hadron Collider [2] is so demanding in terms of both computing power and mass storage that data and CPU's can not be concentrated in a single site and will be geographically distributed (in some cases, data will even be replicated in more than one site) on a computational Grid [3] according to a "multi-tier" model first conceived within the Monarc Project [4]. Moreover, LHC experiments' communities are made of several thousands of people from a few hundreds of institutes spread out all over the world. These people, according to their collaborations on specific physics analysis topics, can constitute highly dynamic Virtual Organizations quite rapidly changing as a function of both time and topology.

---

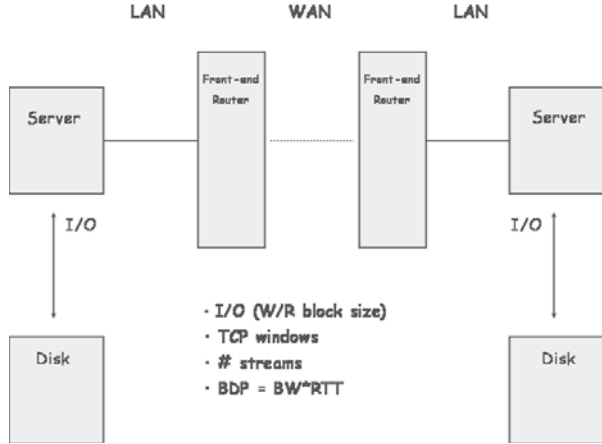* Corresponding author. Email: roberto.barbera@ct.infn.it.

Figure 1. File transfer schema.

For all of the above, the impact of future experiments on Wide Area Networks (WAN) will be absolutely non negligible especially for what concerns the capillarity of bandwidths (down to the "last mile"), quality of service, adaptivity and configurability.

In this paper we report on a series of multi-site data transfer tests performed within the ALICE Experiment [5] on a wide area network test-bed in order to spot possible bottlenecks and pin down critical elements and parameters of actual research networks.

In order to make the tests as realistic as possible, reflecting the real use cases foreseen in the next future, we have taken into account all the aspects of the elements involved in the transfer of a file (the atomic data set considered in this study) and sketched in Fig. 1:

- the local disk Input/Output (I/O) performance;
- the I/O block size;
- the TCP parameters and number of parallel streams;
- the Bandwidth Delay Product (BDP) expressed as the product of the Bandwidth (BW) times the Round Trip Time (RTT).

The paper is organized as follows. The configuration of the wide area network test-bed which has been used is explained in Section 2 while its characterization is reported in Section 3. Section 4 is devoted to the description of the results while summary and conclusions are drawn in Section 5.

## 2. Test-bed lay-out and set-up

Table 1 show the composition of the WAN test-bed on which all tests described in this paper have been carried out. All machines run Linux: some of them are commodity PC's and some others are high performances disk servers. As shown in the Table, the sites belonging to the test-bed also have quite different values of the access or guaranteed (BGA) bandwidths to the Internet, spanning from about 8 Mbit/s to 1 Gbit/s. It is worth noting that, in order to make tests as close as possible to the a real use, the machines chosen were not fully allocated for the tests but they were also concurrently running local simulation and/or analysis jobs. All machines are servers and worker nodes of ALICE farms except ccms.ba.infn.it (see Table 1) which belongs to the CMS Experiment [6].

Figure 2 shows the geographic lay-out of the test-bed. Most of the sites are located in Italy and are connected to the GARR Consortium [7] network, the Italian National Research and Education Network. ALICE sites in Houston (USA) and Prague (Czech Republic) were also included in the test-bed in order to quantitatively study the effectiveness of data transfers as a function of both the number of parallel TCP streams and the TCP parameters in cases of largely different Bandwidth Delay Products (see Section 4.3).

| Site | Max bw/BGA (Mbit/s) | Machines | Exp. |
|---|---|---|---|
| Bari | 28/16 | alicegrid1.ba.infn.it alicegrid2.ba.infn.it ccms.ba.infn.it | ALICE CMS |
| Bologna | 32 | boalice1.bo.infn.it boalice2.bo.infn.it boalice3.bo.infn.it boalice6.bo.infn.it boalice8.bo.infn.it boalice10.bo.infn.it | ALICE |
| Cagliari | 8 | antani.ca.infn.it server3.ca.infn.it | ALICE |
| Catania | 34 | aliserv10.ct.infn.it alipc6.ct.infn.it | ALICE |
| CNAF | 1024 | dell26.cnaf.infn.it wn-04-25-a.cnaf.infn.it | ALICE |
| Houston | 100 | psycho.hpcc.uh.edu | ALICE |
| Padova | 155 | pcalice15.pd.infn.it pcalice16.pd.infn.it pcalice17.pd.infn.it pcalice18.pd.infn.it pcalice19.pd.infn.it pcalice20.pd.infn.it | ALICE |
| Prague | 1000 | golias.farm.particle.cz | ALICE |
| Torino | 155/70 | alifarm01.to.infn.it alifarm02.to.infn.it | ALICE |
| Trieste | 16 | alifarm.ts.infn.it | ALICE |

Table 1. Test-bed composition and maximum network bandwidths.

All file transfers performed in the tests reported in this paper were carried out using bbFTP [8], a very efficient file transfer tool developed at the IN2P3 Computing Centre in Lyon. Among all file transfer tools available, bbFTP has been chosen for the following reasons:

- it is very easy to install and configure;
- it implements its own protocol expressly designed for transfers of very large files (up to 2 GBytes) like those generated by ALICE simulation jobs;
- it can perform multi-stream file transfers;
- different TCP windows and disk I/O block sizes can be easily selected;
- the client wakes the server daemon up only when a file transfer occurs and kills it afterwards so no stray processes are needed to be running on the machines;
- peer-to-peer authentication can be performed both via SSH and X.509 digital certificates;
- file transfers can be set to occur through user-customizable port ranges so the potential impact on strict firewall rules at different sites is reduced to a minimum.
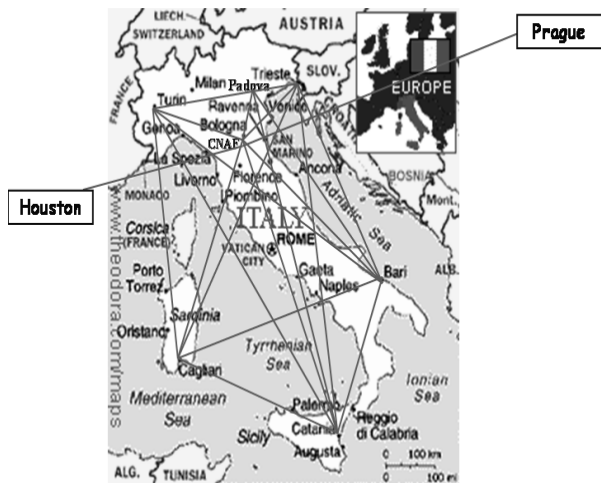


Figure 2. Geographical view of the test-bed lay-out.

SSH private and public keys were created for all machines in the test-bed and public keys were exchanged among all machines in order to let bbFTP working in an automatic, unattended, way without the need of typing passwords to start each file transfer. This also allowed to distribute/update the scripts triggering the file transfers in a centralized way (from CNAF [9]) without any interventions from other system managers, greatly reducing the response time of the different sites which already cover several time zones.

## 3. Test-bed characterization and tools

Very often the nominal bandwidth of a given network route is not just what one can actually use in transferring real files from local disks where data are stored to disks local to the destination node, even in the ideal case of a dedicated point-to-point connection. There are many possible reasons for this discrepancy (Fig. 1): low disk I/O performances, hardware inadequacy or configuration problems of routers or other network elements, TPC/IP protocol inefficiency for high RTT routes, etc. Several tools exist to measure all of these possible sources of inefficiency. In the study reported in this paper, we focused our attention on disk I/O and effective bandwidth measurements. For the first type of measurements we used Bonnie++ [10] and IOzone [11], while for the second one we used Iperf [12] and Netperf [13]. Standard configuration of both TCP stack and disk I/O parameters in Linux are kept constant to the default values in this phase.

| Machine | Write/Read (MBytes/s) | |
|---|---|---|
| | Bonnie++ | IOzone |
| boalice8.bo.infn.it | 5/3 | 5/5 |
| server3.ca.infn.it | 43/32 | 45/61 |
| aliserv10.ct.infn.it | 57/25 | 27/34 |
| pcalice19.pd.infn.it | 5/5 | - |
| alifarm02.to.infn.it | 31/53 | 40/95 |
| alifarm.ts.infn.it | 27/34 | 28/36 |

Table 2. Input/output measurements on a representative subset of the test-bed nodes performed with Bonnie++ and IOzone.

IOzone and Bonnie++ are filesystem benchmark tools. They generate and measure a variety of file operations such as write/read of characters, write/read of blocks of memory, random write/read, etc. In order to get reliable measurements, tests were performed creating files having a total size of at least twice the total RAM memory available on the machine under consideration. Results are summarized in Table 2. Keeping in mind that the measurements were done on non-dedicated machines, the values of disk access speed obtained with the two different tools are quite consistent and, more importantly, they show how disk access is not expected to be a bottleneck for the file transfers over WAN.

| Machine | BW 1 | BW 2 | BW 4 | BW 8 | BW 16 | BW 32 |
|---|---|---|---|---|---|---|
| boalice8 | 76 | 77 | 79 | 84 | 86 | 87 |
| server3 | 12 | 21 | 22 | 21 | 21 | 22 |
| aliserv10 | 9 | 15 | 18 | 18 | 19 | 20 |
| pcalice19 | 26 | 51 | 87 | 92 | 93 | 94 |
| alifarm02 | 27 | 50 | 57 | 61 | 64 | 69 |
| alifarm | 14 | 18 | 18 | 18 | 19 | 19 |

Table 3. Bandwidth measurements performed with Iperf on a representative subset of test-bed machines.

Iperf and Netperf are tools to measures network performances both for TCP and UDP protocols. In particular, Iperf is quite versatile and, besides simple bandwidth measurements, it is able to perform a large variety of functions such as TCP tuning and multi-streaming. We performed a series of measurements of the actual maximum capacity of the network in a "star" configuration, i.e. using one of the CNAF machines as Iperf/Netperf server, since the CNAF link to the Internet has the largest bandwidth among the sites belonging to the test-bed (see Table 1).

| Machine | BW 1 | BW 2 | BW 4 | BW 8 | BW 16 | BW 32 |
|---|---|---|---|---|---|---|
| boalice8 | 30 | 44 | 65 | 80 | 81 | 86 |
| server3 | 13 | 18 | 22 | 22 | 22 | 23 |
| aliserv10 | 9 | 16 | 19 | 20 | 22 | 22 |
| pcalice19 | 26 | 51 | 87 | 92 | 93 | 97 |
| alifarm02 | 28 | 41 | 46 | 55 | 61 | 65 |
| alifarm | 14 | 17 | 18 | 18 | 17 | 19 |

Table 4. Bandwidth measurements performed with Netperf on a representative subset of test-bed machines.

Results are summarized in Table 3 for Iperf and in Table 4 for Netperf, respectively. They are similar, especially when a large number of parallel TCP streams was used. In fact, the BW1, BW2, …, BW32 symbols refer to the bandwidth values (expressed in Mbit/s) obtained with 1, 2, …, 32 parallel streams.

## 4. Results

Tests have been performed with two different topologies of the WAN test-bed: flat and multi-tier. The flat topology refers to the chaotic case of really distributed analysis where files of different sizes are moved among all the sites of the test-bed. The multi-tier topology, on the other hand, describes the more "ordered" case where productions of simulated events are executed at some Tier-1/2 sites and then these events are remotely analyzed by users located at Tier-3/4 sites. Results of tests performed with flat and multi-tier topology are reported in Section 4.1 and 4.2, respectively.

### 4.1 Flat topology

The flat topology has been used in the first phase of testing and it is aimed to find the possible bottleneck and general limitations when one tries to use the network to move large amounts of data chaotically, i.e. like in a real distributed analysis on a Grid.

The traffic is generated through bbFTP with a customizable number of parallel streams (1, 2, 4, 6, 8, 16, and 32 parallel stream can be chosen). Data consist of a set of files of three different sizes (300 MBytes, 800 Mbytes, and 1.6 GBytes) that are randomly chosen for each transfer. Between consecutive file transfers, a time delay can be randomly set between 0 and a customizable maximum (1 minute and 5 minutes have been used as maxima).
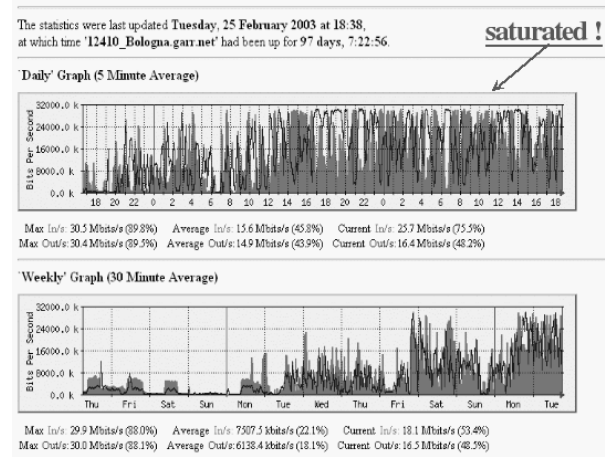


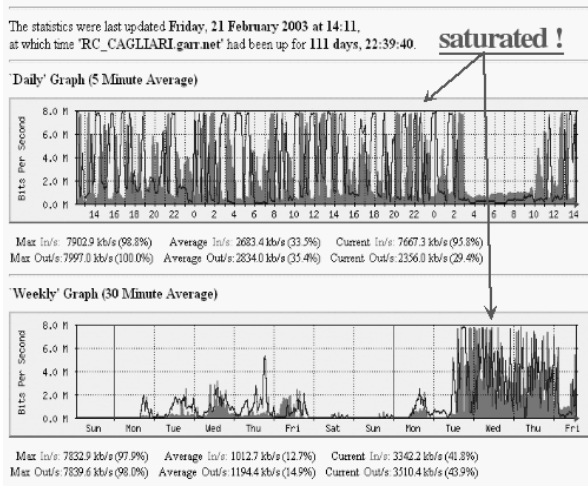Figure 3. Traffic through the INFN Bologna router.

The statistics were last updated Friday, 21 February 2003 at 14:11, at which time 'RC_CAGLIARI.garr.net' had been up for 111 days, 22:39:40.

**saturated !**

'Daily' Graph (5 Minute Average)

Max In/s: 7902.9 kb/s (98.8%)  Average In/s: 2683.4 kb/s (33.5%)  Current In/s: 7667.3 kb/s (95.8%)
Max Out/s: 7997.0 kb/s (100.0%)  Average Out/s: 2834.0 kb/s (35.4%)  Current Out/s: 2356.0 kb/s (29.4%)

'Weekly' Graph (30 Minute Average)

Max In/s: 7832.9 kb/s (97.9%)  Average In/s: 1012.7 kb/s (12.7%)  Current In/s: 3342.2 kb/s (41.8%)
Max Out/s: 7839.6 kb/s (98.0%)  Average Out/s: 1194.4 kb/s (14.9%)  Current Out/s: 3510.4 kb/s (43.9%)

Figure 4. Traffic through the INFN Cagliari router.

Each machine in the test-bed alternated downloads and uploads and every time the source or destination machine (partner in that specific file transfer) was randomly chosen with a probability weighted over the maximum bandwidth of the site that machine belongs to.

Typical figures of the network occupancy for some sites belonging to the test-bed are shown in Figs. 3, 4, and 5. Data come from the official site of the GARR Network Operation Center (NOC) [14] and are relative to file transfers with 8 TCP parallel streams. The file transfer tests spotted immediately the severe limitations of bandwidth, especially taking into account that only a tiny fraction of the machines/servers already installed or foreseen to be installed in the next future at the different ALICE sites in Italy were used to set-up the test-bed. It must also be said, however, that, although tests lasted for several days, absolutely no service interruption occurred in any of the test-bed sites.



The statistics were last updated Tuesday, 18 February 2003 at 19:42, at which time 'RC_CATANIA.garr.net' had been up for 141 days, 6:38:55.

**heavy traffic !**

'Daily' Graph (5 Minute Average)

Max In/s: 30.0 Mb/s (97.6%)  Average In/s: 14.6 Mb/s (47.4%)  Current In/s: 1763.4 kb/s (5.7%)
Max Out/s: 29.0 Mb/s (94.5%)  Average Out/s: 13.8 Mb/s (45.1%)  Current Out/s: 1550.3 kb/s (5.0%)

'Weekly' Graph (30 Minute Average)

Max In/s: 26.4 Mb/s (85.9%)  Average In/s: 8280.0 kb/s (27.0%)  Current In/s: 8557.7 kb/s (27.9%)
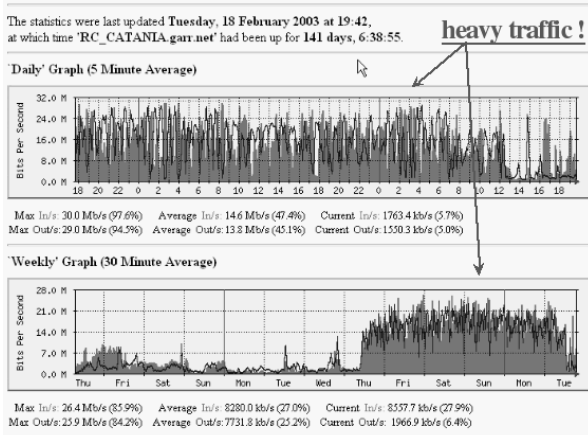Max Out/s: 25.9 Mb/s (84.2%)  Average Out/s: 7731.8 kb/s (25.2%)  Current Out/s: 1966.9 kb/s (6.4%)

Figure 5. Traffic through the INFN Catania router.

Just after this first phase of testing with flat topology the GARR Consortium provided a quick and big increase in bandwidth in the test-bed sites which had shown a network capacity clearly inadequate. The main upgrades concerned INFN-Bologna with an increase of the maximum speed of data transmission from 32 Mbit/s to 100 Mbit/s, INFN-Cagliari from 8 Mbit/s to 32 Mbit/s, and INFN-Trieste from 16 Mbit/s to 24 Mbit/s.

## 4.2 Multi-tier topology

The second phase of tests consisted in a simulation of the network traffic coming from a distributed production of data and the following analysis in the context of a multi-tier computing model as the one envisaged within the Monarc Project.

The use case for these multi-tier tests was the production of 5000 central Pb+Pb simulated events at LHC energy which are expected to produce a total output of about 9 TBytes. We made the hypothesis to produce 60% of these data at the ALICE (and INFN) Tier-1 Centre located at CNAF and 20% at the two ALICE Tier-2 Centres located at Catania and Torino. We also realistically assumed to replicate at the Tier-1 Centre the data produced at the Tier-2 Centres in order to store there the full production set. The first step is then the upload of the Tier-2 data into the Tier-1 storage facility, with a granularity of 1.8 GBytes files (see Fig. 6).
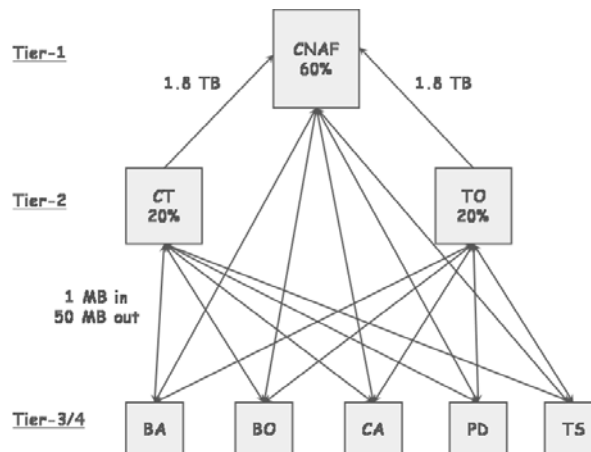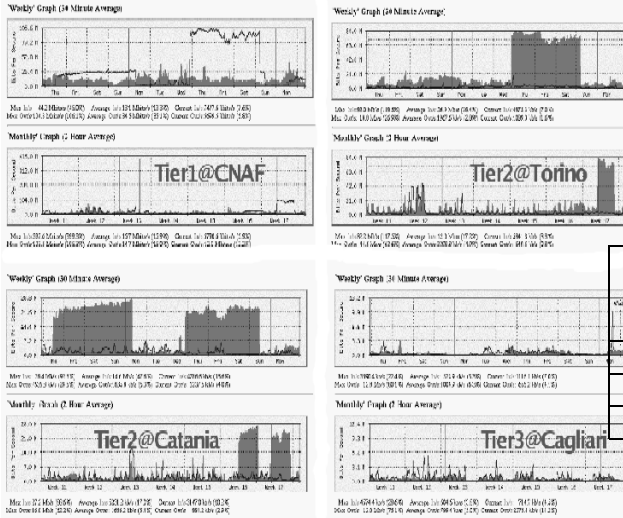


Figure 6. Multi-tier test-bed lay-out.

Figure 7. Traffic through some of the Tier Centres' routers.

trip time (RTT), it represents a quantitative indication of how much data can stand in a given network path across two remote sites. The standard Linux TCP configuration is not optimized for transferring files along a large BDP route and that is why we performed a series of tests changing the TCP windows in order to maximize the throughput.

| Site | RTT (ms) from CNAF | BW (Mbit/s) from CNAF | BDP (MBytes) |
|------|------|------|------|
| Catania | 25 | 25 | 0.008 |
| Houston | 140 | 70 | 1.2 |
| Prague | 20 | 25 | 0.6 |

Table 5. Some characteristics of the three network routes we considered in order to study the effect of TCP windows, number of streams, and I/O block size on throughput.

When a certain fraction (we set this value at 80%) of Tier-2 data have already been replicated at Tier-1, we assumed the Tier-3's start the analysis sending jobs both to Tier-1 and Tier-2's. Since we were interested in the network impact of this activity, analysis jobs have been simulated just by time intervals between an upload of a small input file (1 MB), which would contain the input instructions for the analysis, and the download of the results, a 50 MB file that could contains histograms or even more complicated data structures. The duration of the jobs was made to depend on where the job is sent: 5 minutes if it runs at a Tier-2 centre, 15 minutes if it runs at the Tier-1.

The whole schema of multi-tier test, as it has been described above, is reported in Fig. 6 and some traffic statistics from GARR NOC are shown in Fig. 7. Results confirm the expectation that, within a multi-tier model, Tier-3's experience quite a low traffic and, on the contrary, Tier-1 and Tier-2's have to sustain a really heavy traffic, and can easily go to saturation.

## 4.3 TCP streams and parameters dependency

All tests described in Section 4.1 and 4.2 have been carried out with default values of the parameters of the TCP stack and with a fixed number (8) of parallel TCP streams. In this Section we will describe the third and last phase of our test concerning throughput measurements for really wide area file transfers between CNAF and Houston and between CNAF and Prague with different TCP customizations. Some of the characteristics of these two routes are reported in Table 5. Catania has been included as a low Bandwidth Delay Product (BDP) reference site. Since BDP is defined as the product of bandwidth (BW) and round

The base line of all third-phase tests was the use of all meaningful parameters (number of streams, I/O block size, TCP windows) in order to determine which conditions enable the maximum throughput.

On this purpose, we carried out bbFTP transfers of files with different sizes, spanning from 10 MBytes to 800 MBytes, and measured the throughput as a function of both the number of treams and the TCP windows (to simplify the analysis of the results, we kept receiver and sender windows equal). In order to take into account throughput fluctuations, due for example to concurrent background traffic or to the routing, we repeated transfers many times (at least five) and, for a given set of parameters, we assumed as final result the average with its statistical error.

Figures 8, 9, and 10 show the results of these throughput measurements for the route CNAF-Catania as a function of the TCP windows size and for different numbers of parallel TCP streams. Saturation comes already with single-stream transfers for large files over a TCP windows threshold of about 128 kBytes. Rising the number of parallel TCP streams allows to maximize the throughput also for small size files producing the effect of moving the TCP window threshold towards smaller values. For all these measurements I/O block size has been kept constant to a value of 256 kBytes.
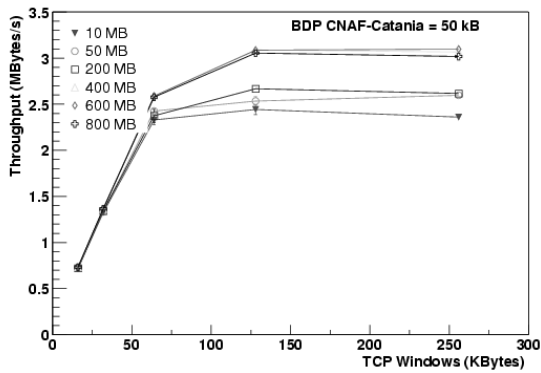
Figure 8. Throughput measurements for the route CNAF-Catania as a function of the TCP windows size and with a single TCP stream.
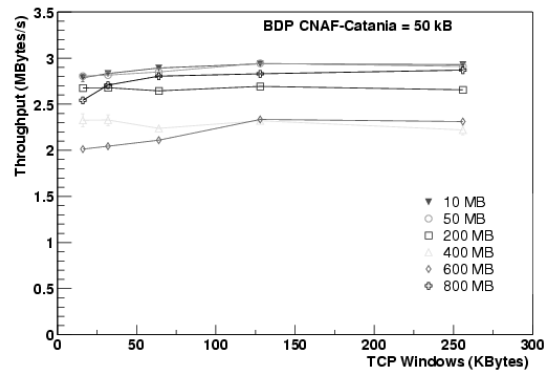


Figure 10. Throughput measurements for the route CNAF-Catania as a function of the TCP windows size and with four parallel TCP streams.
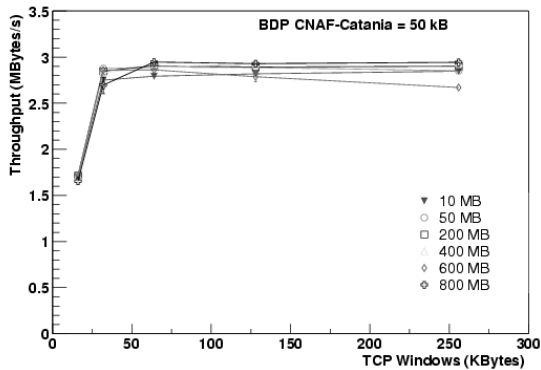


Figure 9. Throughput measurements for the route CNAF-Catania as a function of the TCP windows size and with two parallel TCP streams.
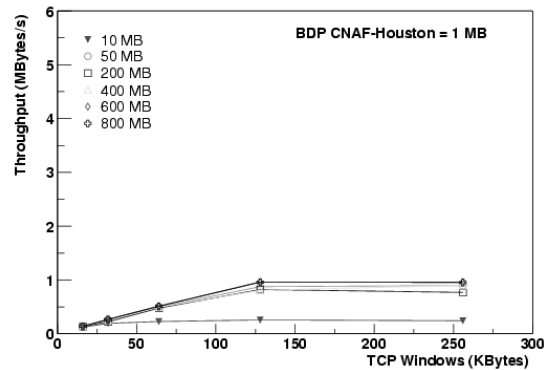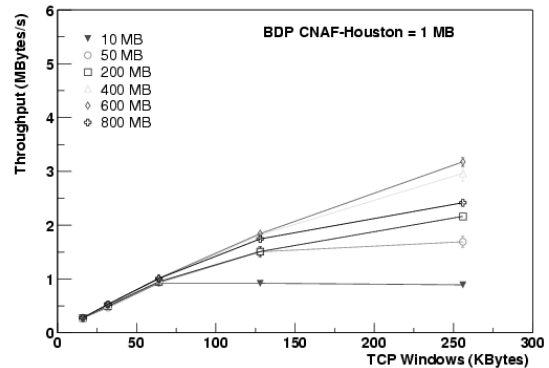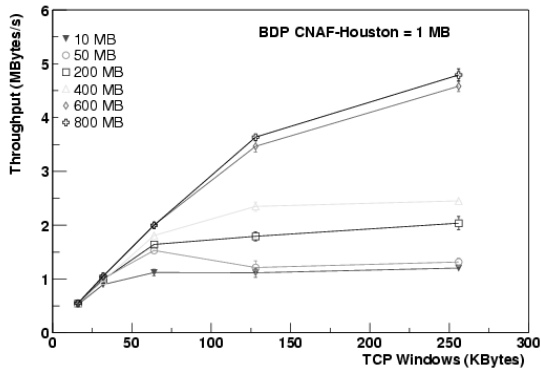


Figure 11. Throughput measurements for the route CNAF-Houston as a function of the TCP windows size and with a single TCP stream.

Figures 11, 12, 13, and 14 show the results of the same measurements for the route CNAF-Houston. In this case there is a quite large RTT (and then BDP) and the effect of modifying TCP windows is more evident. Indeed, increasing TCP windows from smaller to larger values, data rate rises by a factor of 5. Saturation is obtained again for a TCP window threshold of about 128 kBytes but it is not reached for small files and for larger files it comes only with a higher number of streams (six instead of just one). I/O block size was 256 kBytes also in this case.



Figure 12. Throughput measurements for the route CNAF-Houston as a function of the TCP windows size and with two parallel TCP streams.

Figure 13. Throughput measurements for the route CNAF-Houston as a function of the TCP windows size and with four parallel TCP streams.
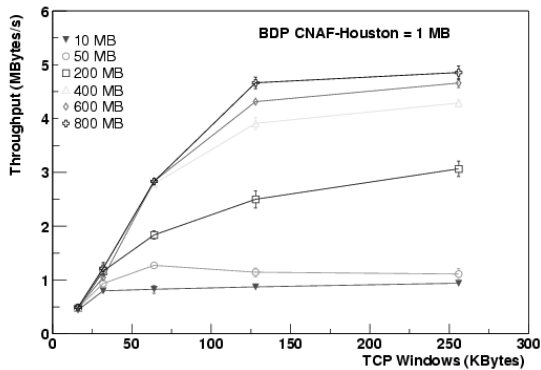


Figure 14. Throughput measurements for the route CNAF-Houston as a function of the TCP windows size and with six parallel TCP streams.

The last route we considered was CNAF-Prague. It this case the RTT is smaller than the CNAF-Houston case but the BDP is still rather large, thanks to the fact that on both sides of the connection there are Gigabit Ethernet interfaces. Again, the increase of throughput caused by the use of large TCP windows is evident (see Figs. 15, 16, 17, and 18), and rising the number of streams gives good performances also for small files. These results were also obtained with a I/O block size of 256 kBytes.

An additional set of measurements were made for the route CNAF-Prague and their results are shown in Figs. 19 and 20. In this last case, the I/O block size was increased from 256 kBytes to 1 MByte and the default Linux TCP parameters were modified according to Ref. [15]. Figure 19 shows the results obtained with a single TCP stream while Fig. 20 shows the results obtained with 4 parallel TCP streams. In both cases the changes allow for a larger throughput

with respect to the cases where the default values of the I/O block size and TCP parameters were used.
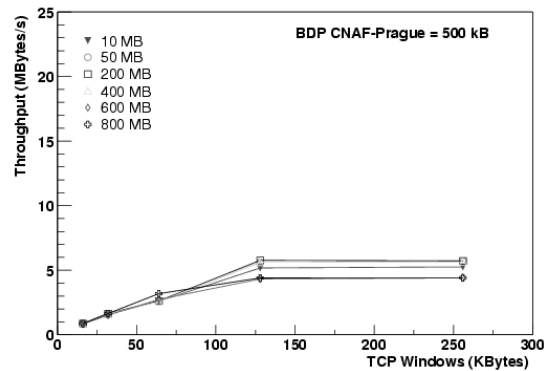


Figure 15. Throughput measurements for the route CNAF-Prague as a function of the TCP windows size and with a single TCP stream.
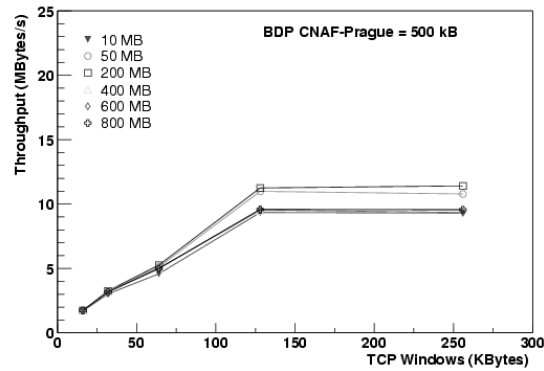


Figure 16. Throughput measurements for the route CNAF-Prague as a function of the TCP windows size and with two parallel TCP streams.
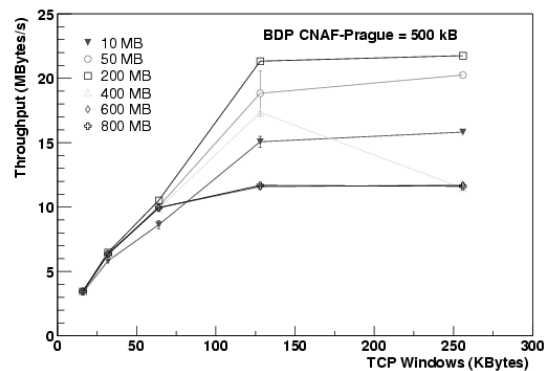


Figure 17. Throughput measurements for the route CNAF-Prague as a function of the TCP windows size and with four parallel TCP streams.
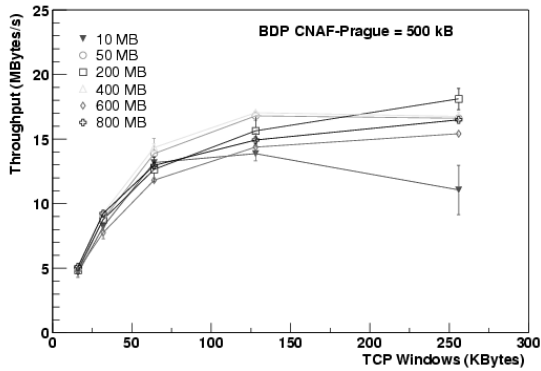
Figure 18. Throughput measurements for the route CNAF-Prague as a function of the TCP windows size and with six parallel TCP streams.
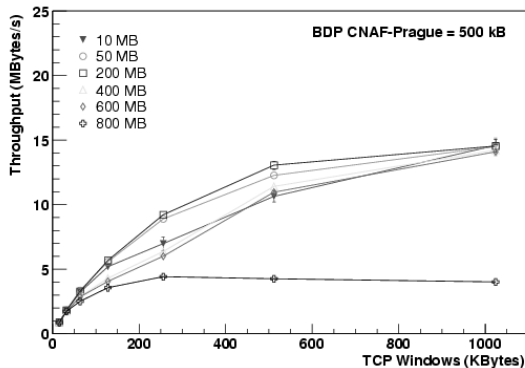


Figure 19. Throughput measurements for the route CNAF-Prague as a function of the TCP windows size and with a single TCP stream. The I/O block size was set to 1 Mbyte and TCP parameters were modified according to Ref. [15].
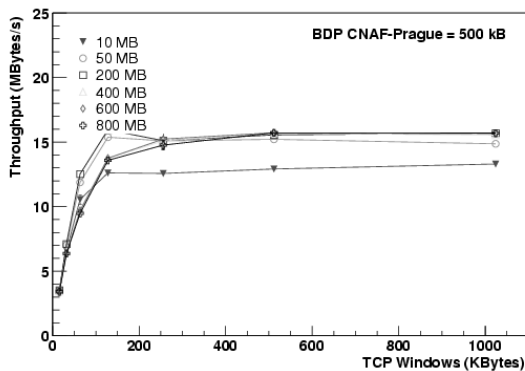


Figure 20. . Throughput measurements for the route CNAF-Prague as a function of the TCP windows size and with four parallel TCP streams. The I/O block size was set to 1 Mbyte and TCP parameters were modified according to Ref. [15].

## 5. Summary and conclusions

The computing model of next generation high energy physics experiments planned at the CERN LHC is inherently "distributed" and will most be based implemented on computational Grids. Thus, both coordinated and chaotic analysis activities will have a big impact on networks and will strongly rely on their robustness and reliability.

In this paper we have reported the results of realistic wide area network tests performed on a large scale test-bed and aimed at spotting possible bottlenecks and pin down critical elements and parameters of actual research networks.

Conclusions can be summarized as follows:

- "network" not only means the large bandwidth of international links but also, and more important, reliable end-to-(many)ends connections (the so-called "last mile problem" should be addressed and hopefully solved);

- scientific Virtual Organizations are very dynamical as a function of both space and time (they will be driven by physics topics) so best effort and over-provisioning could not always be the best solutions; quality of services, bandwidth-on-demand and advanced reservation will be key issues of future networks;

- the analysis activity, requiring the sum of large fraction of the total data sample, can not be based on the transfer of all the distributed input to a single site which executes the algorithm, as that approach would saturate network connections, no matter their bandwidth; application specific code must interact with Grid middle-ware services in such a way that parallel analysis of sub-samples, grouped according to the input geographical distribution, followed by output merging be possible.

## 6. Acknowledgments

## 7. References

[1] http://www.cern.ch

[2] http://www.cern.ch/lhc

[3] Foster I, Kesselman C, editors. The GRID: blueprint for a new computing infrastructure. San Francisco: Morgan Kaufmann; 1999

[4] http://www.cern.ch/MONARC

[5] http://alice.web.cern.ch/Alice/AliceNew

[6] http://cmsinfo.cern.ch/Welcome.html

[7] http://www.garr.it/

[8] http://doc.in2p3.fr/bbftp

[9] http://www.cnaf.infn.it

[10] http://www.coker.com.au/bonnie++

[11] http://www.iozone.org

[12] http://dast.nlanr.net/Projects/Iperf/

[13] http://www.netperf.org/netperf/NetperfPage.html

[14] http://www.noc.garr.it/

[15] http://www-didc.lbl.gov/TCP-tuning