

iGrid2002 Demonstration: Bandwidth from the Low Lands

R. Les. Cottrell, Antony Antony, Connie Logg, and Jiri Navratil

Abstract—We report on a demonstration of several complementary high performance end-to-end active network throughput measurement tools. We also demonstrate sending high-speed data from 4 hosts at iGrid2002 to over 30 hosts in 10 countries to simulate a high energy physics experiment distributing data to collaborators. The demonstration utilized the high-speed, long latency, trans-Atlantic network set up for iGrid2002 in Amsterdam during September 2002.

Index Terms—iGrid2002, high-throughput, measurement tools, monitoring, networks, tuning, TCP.

I. INTRODUCTION

The avalanche of data already being generated by and for new and future High Energy and Nuclear Physics (HENP) experiments demands a new strategy for how data is collected, shared, analyzed and presented. For example the SLAC BaBar experiment [1] and JLab [2] are each already collecting over a Tbyte/day, and BaBar expects to increase by a factor of 2 in the coming year. The SLAC BaBar and Fermilab CDF [3] and D0 [4] experiments have already gathered well over a Petabyte of data, and the LHC experiments [5] expect to collect over ten million Tbytes. The strategy being adopted to analyze and store this unprecedented amount of data is the coordinated deployment of Grid technologies such as those being developed by the European Data Grid [6], Particle Physics Data Grid [7] and the Grid Physics Network [8]. It is anticipated that these technologies will be deployed at hundreds of institutes that will be able to search out and analyze information from an interconnected worldwide grid of tens of thousands of computers and storage devices. This in turn will require the ability to sustain over long periods the transfer of large amounts of data between collaborating sites with relatively low latency.

The opportunity promised by access to the high speed, long latency, trans-Atlantic network put together specifically for iGrid2002 [9] in Amsterdam during September 2002,

motivated us to demonstrate high throughput measurements by a variety of state-of-the art tools, across this network.

Our iGrid2002 project/demonstrations were designed to show: the current data transfer capabilities from iGrid2002 to over 30 HENP, grid or network monitoring sites with high performance links, worldwide; to compare the performances from various monitoring sites; and to demonstrate and compare light and heavyweight methods of measuring performance. Further we deliberately choose to do make these measurements using standard network (TCP/IP) implementations, standard Maximum Transfer Units (MTUs) of 1500 bytes, and with no efforts to try and provide preferred quality of service. In a sense the site at iGrid2002 was acting like a HENP tier 0 or tier 1 site [10] (an accelerator or major computation site) in distributing copies of the raw data to multiple replica sites.

In this paper we first describe the configurations of the measuring equipment, the network and the remote hosts that were setup for this demonstration. Then we describe the various demonstration tools, how they were set up and show example screen-shots of the visualizations. We follow this with a discussion of the results obtained and close with conclusions.

II. SETUP

A. Hardware & Networking

We had 5 Linux hosts, each with a 64bit 66MHz PCI bus, a Fast Ethernet (100 Mbits/s) connection and one or two 1GE (Gigabit Ethernet) Network Interface Cards (NICs). More details of the host configurations are given in Table 1. GE NIC I was used for all measurements. We did not use GE NIC II in any host during iGrid2002. The hosts were connected to a Cisco 6509 switch. The 6509 was connected at 10Gbits/s to the SURFnet core. The major external connections from the SURFnet core were at 10 Gbits/s and 2.5Gbits/s to StarLight, 2.5Gbits/s to GEANT, 2.5 Gbits/s to CERN, 622Mbps to StarLight for U.S. research networks (ESnet and Abilene), 1 Gbits/s to StarTAP, 2 x 1 Gbits/s to the Internet at large (mainly .com traffic) and 1 Gbits/s peering across the Amsterdam Internet Exchange. We ensured that the various TCP ports needed by our applications were not blocked at the iGrid2002 border.

Manuscript submitted October 31, 2002. This work was supported in part by the Director, Office of Science, Office of Advanced Scientific Computing Research, Mathematical, Information, and Computational Sciences Division under the U.S. Department of Energy. The SLAC work is under Contract No. DE-AC03-76SF00515..

Les Cottrell, Connie Logg and Jiri Navratil are with the Stanford Linear Accelerator Center, 2575 Sand Hill Road, Menlo Park, CA 94025 (emails: cottrell@slac.stanford.edu, cal@slac.stanford.edu, jiri@slac.stanford.edu).

Antony Antony is with the national Institute for Nuclear and High Energy Physics (NIKHEF), Kruislaan 409, 1009 SJ Amsterdam, Netherlands (email: antony@nikhef.nl), funded by the IST Program of the European Union (grant IST-2001-32459) via the DataTAG project.

Host name	CPU # x GHz	RAM GB	GE NIC I	Linux Kernel	GE NIC II
keeshond	2 x 1	2	SK9843	2.4.19 net100	3c985
stier	2x 2	2	Intel EE fiber	2.4.19 net100	3c996
haan	1 x 0.8	0.5	3c985	2.4.19 net100	
hp3	2 x 2.4	1	Intel EE Copper	2.4.18	
hp4	2 x 2.4	1	Intel EE copper	2.4.19 net100	

Table 1. Host configurations

We made contact with system administrators or network people at other sites who provided us with access to one or more hosts at their sites. The list of contacts and sites can be found in [11]. These sites were distributed over 10 countries and there were 34 *remote hosts* that we set out to measure performance to. The routes to the sites and the average throughputs in Mbits/s achieved from iGrid2002 demonstration are shown in Fig. 1. The names of remote hosts that are using 100Mbit/s NICs are written in italics, the others had 1Gbit/s NICs. The bold face numbers on the links between Internet Service Provider (ISPs) clouds such as GEANT and SurfNet indicate the speed of these links in Gbits/s. More details on the network connections to iGrid2002 can be found in [12]. All hosts were setup to enable the use of large (≥ 1 Mbyte) TCP windows/buffers. The monitoring hosts were set up to flush the Linux TCP buffer caches.

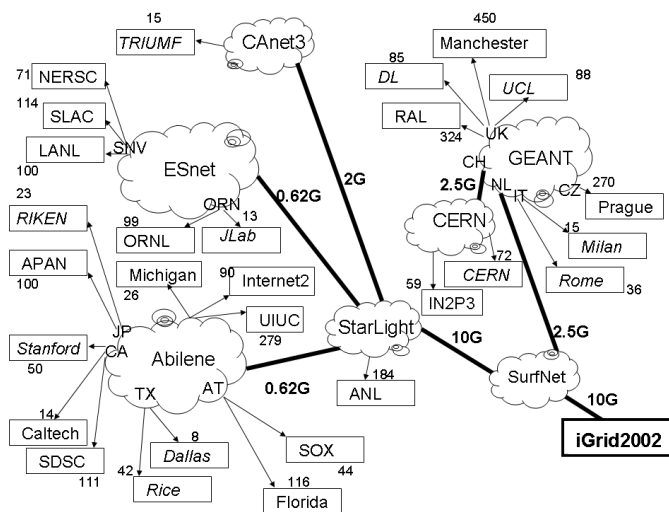


Figure 1: Routes and Mbits/s achieved from iGrid2002 to the remote sites.

B. Software & accounts

Each monitoring host in our demonstration was installed with the software described in [13]. To simplify managing the

software on multiple hosts we used NFS to mount the directories that were stored centrally on keeshond. The NFS connections utilized the Fast Ethernet connections on all hosts except stier that mounted over its GE NIC I.

III. DEMONSTRATIONS

There have been many tools developed for network performance measurements in the last few years [14], [15]. For this demonstration we were concerned only with “active” measurement tools, i.e. tools that inject traffic/probes into the network to make the measurements. Most current active monitoring is based on simple round trip (e.g. ping [16], PingER [17], AMP [18]) or one-way delay measurements (e.g. Surveyor [19], RIPE [20]); route discovery tools such as traceroute [21]; packet pair techniques such as pathload [22], pipechar and netest2 [23]; or injecting streams of UDP or TCP traffic (e.g. iperf [24]) into the network. Increasingly user applications for data transfer with known file sizes are also being utilized. Such applications include bbcp [25], bbftp [26] and GridFTP [27]. From first to last in the above list, each of these tools typically injects increasingly large amounts of network traffic, and measures increasingly closer to what the end user would expect for an end-to-end application. There is no single tool that would measure all the network performance metrics, in fact there are multiple ways to measure a single metric. For example, when measuring throughput one can use Round Trip Times (RTT) and losses together with the Mathis formula [28] to deduce TCP bandwidth, or one can use packet pair techniques, or iperf, or an application. Thus, in this demonstration, we used several tools that we have developed, and which are described in the following sections.

A. PingWorld

We demonstrated the PingWorld [29] Java applet to illustrate RTT and losses from iGrid2002 to over 30 different countries. PingWorld uses the ping facility in the client to send a default sized ping once a second to remote hosts. These remote hosts were chosen from experiences with the PingER project measurements. The requirements for the remote hosts were that they were available, and representative of regions of the world or major national research networks. The results, seen in Fig. 2, are rendered by PingWorld as multiple time-series, each representing the round trip delay to the remote host. Adjacent points are joined by lines and broken lines indicate packet loss. The time-series are gathered into groups of time-series plots, one for each of 8 regions of the world. At the bottom left a table indicates the most recent RTT, updated once a second, for each host, the color of the square indicating the RTT (red for long, blue for short etc.)



Figure 2: PingWorld demonstration of RTT to remote hosts around the world.

B. Tomography

To illustrate the routes from a measurement host to the remote hosts, we developed a route (topology) and performance (tomography) visualization tool. The tool makes traceroutes to the defined remote hosts and analyzes the traceroute outputs to create an internal database for the routes with common paths/routers identified. From this database a map of the routes is created. A simple illustrative example of the routes from iGrid2002 is shown Fig 3. The top node (circle) identifies the iGrid2002 monitoring host. By default, the colors of the nodes and edges in Fig 3 are colored by Internet Service Provider (ISP). The user can select to color the nodes and edges by RTT or by bandwidth (measured between the end points). The user can also select a node to drill down to more detailed information about the node. To facilitate viewing complex maps, the user can select to see subtrees and individual routes with the full node names displayed.

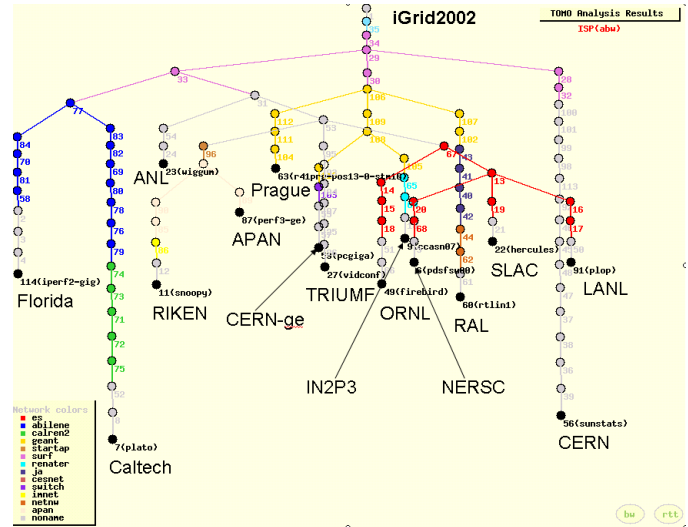


Figure 3: Routes from iGrid2002 visualized by the topology tool.

C. Bandwidth Estimation

To demonstrate the available bandwidth in real-time we needed a tool that was able to run on a monitoring host and quickly make a reasonable estimate of the available bandwidth for 10 to 20 remote hosts in less than 10 seconds. At the same time it needed to do this with limited impact on the network. We evaluated the ability of pipechar, pathrate and pathload to make such measurements. However, we found that current implementations failed for high-speed paths (e.g. over 155Mbits/s) and also pipechar and pathrate could take several minutes to run.

We therefore developed a new available bandwidth estimator (ABWE) for monitoring networks. It is based on packet pair techniques and was designed to work in continuous mode, with high performance paths, to meet our real-time and low network intrusion goals. Similar to other tools of this type, it sends packet pairs with known delays (resolution of 1 μ s) from the client (monitoring host) to the server (remote host). The server measures the inter-packet delays, thus synchronized clocks are not necessary. The number of pairs and the packet size can be selected. For iGrid2002 we chose to use only 20 pairs and 1000byte packets as a reasonable compromise between intrusiveness, test time, adequate statistics and also a reasonable packet size match for links making relatively heavy use of large file transfers.

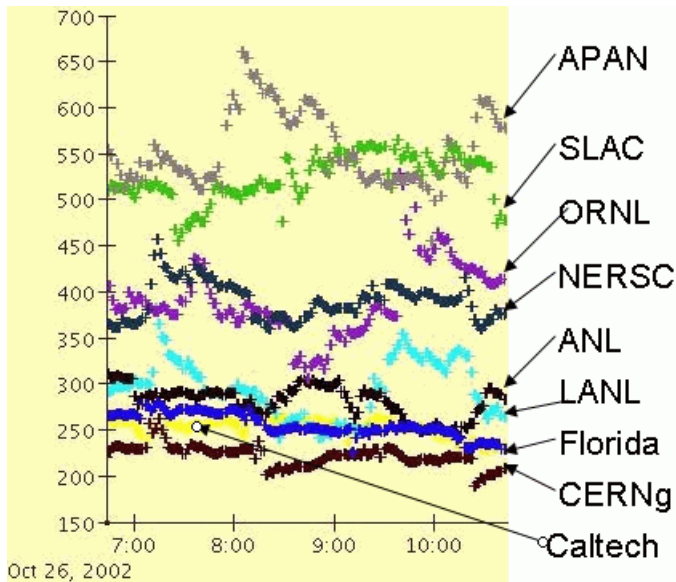


Figure 4: Bandwidth from SLAC to "fast hosts" measured and visualized by ABWE.

For the demonstration we selected 16 remote hosts that had high throughput from iGrid2002 as measured by iperf in TCP mode. These hosts were divided into two groups: the "slow hosts" with throughputs up to 200 Mbits/s and the "fast hosts" with speeds of hundreds of Mbits/s. Real-time time-series bandwidth plots for the two groups were created using the Universal Time History (UTH) package [30]. These were displayed, together with the current packet pair delays, superimposed on a map of the world (similar to the PingWorld demonstration above). Examples of a bandwidth time-series plot for the fast hosts are shown in Fig 4. We did not keep an actual example from iGrid2002, so this example is measured from SLAC. The bottom axis is the time in hours. The points associated with each host were distinguished by colors, and had a colored legend which would not reproduce well in black and white. So we have removed the legend and added labels to help identification.

D. Sequential iperf and Application Measurements

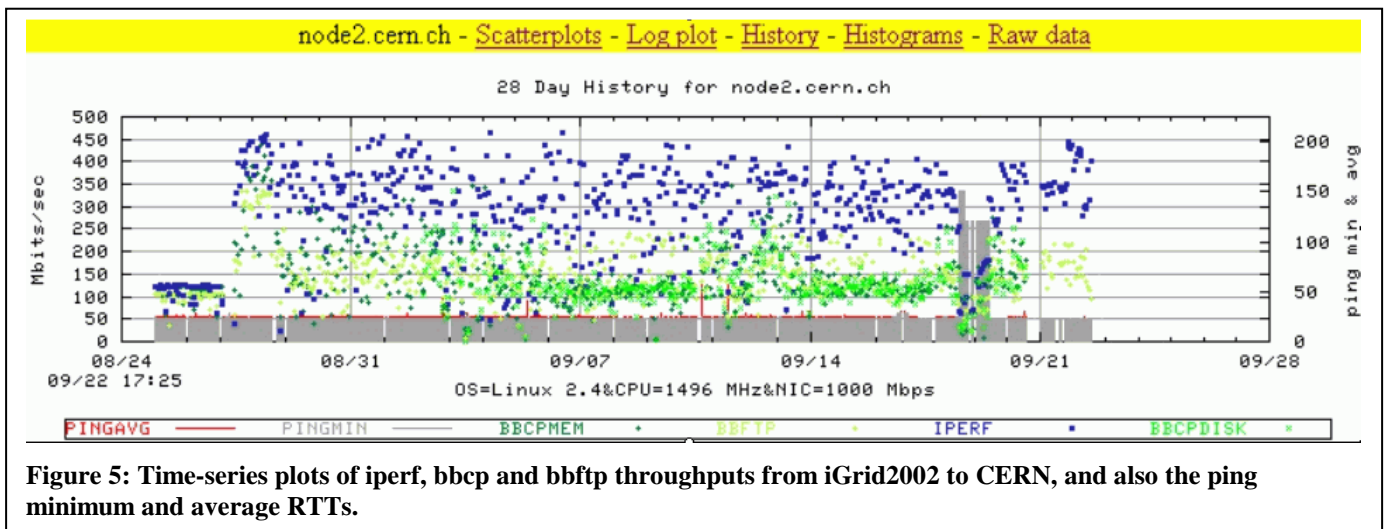


Figure 5: Time-series plots of iperf, bbcp and bbftp throughputs from iGrid2002 to CERN, and also the ping minimum and average RTTs.

To make and demonstrate more accurate (from the end-user viewpoint) measurements of throughput, we deployed the IEPM-BW [31] ssh based active end-to-end performance measurement toolkit on our hosts at iGrid2002. In sequential mode each remote host was monitored once every 90 minutes (a cycle). For each host, we measured the ping response time, the route, the iperf TCP throughput, the bbcp memory-to-memory (/dev/zero to /dev/null) throughput, the bbcp disk-to-disk throughput and the bbftp throughput. Each measurement (ping, iperf, etc.) apart from the traceroute was made for about 10 seconds.

The results from these measurements were recorded in log files. After each cycle, the log files were analyzed to extract the relevant information and this was recorded in tables. The tabular data was then further analyzed to create short and long-term time-series, scatter-plots, histograms, tables of the measurements and other relevant information (e.g. how successful the measurements were). This information was made available via the web with access from a top level page, see for example reference [31]. An example of a time-series plot of the iperf, bbcp memory-to-memory, bbcp disk-to-disk, bbftp and minimum and average RTTs is seen in Fig 5.

The optimal window sizes and number of streams to use for the iperf measurements were determined previously for each remote host by transferring iperf TCP data for 10 seconds to the remote host. The choice of 10 seconds was a reasonable compromise between network impact, time required for each measurement and allowing sufficient time for the transfer to reach reasonable stability (see reference [32] for more details). For each transfer we used a window size selected from 8Kbytes to 4Mbytes, and for each window size we varied the number of parallel data streams from 1 to 120. We then plotted the iperf/TCP reported throughput versus the number of streams for each of the window sizes. An example of such a plot is shown in Fig. 6.

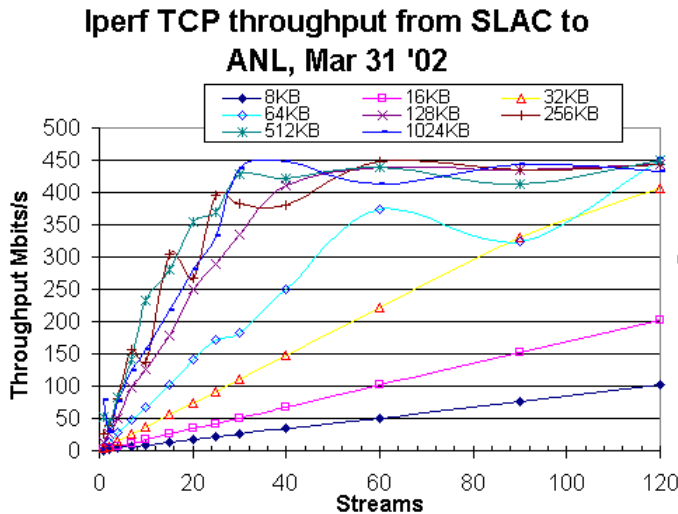


Figure 6: Ten second iperf TCP throughputs from SLAC to ANL as a function of TCP window size and streams.

From such plots, we selected a windows-streams combination that achieved about 80-90% of the maximum throughput measured/achievable, while minimizing the number of streams. We wished to minimize the number of streams since each stream consumes resources in the monitoring and remote hosts. Identical settings of windows and streams were used for iperf, bncp and bbftp for each remote host.

To illustrate the performance of the sequential tests in real-time, we used UTH to read the monitoring hosts bytes transferred (using the Unix `ifconfig` utility) at two second intervals and displayed the in and out bytes/s for the last 120 seconds in a time-series in real-time. In addition in a separate window we displayed the name, and properties (IP address, cpu power, OS etc.) of the current remote host, measurement type, windows and streams, and RTT, (see Fig. 7 for an example). It shows the throughput from the SLAC monitoring host. The bumps in throughput are for the: iperf TCP, bncp memory-to-memory, bncp disk-to-disk and bbftp tools. The duration of the iperf bump is ~ 10 seconds. The bncp memory-to-memory application runs for 15 seconds, so the bump is longer. For bncp disk-to-disk and bbftp the applications terminate when the file is transferred or they are timed out after 15 seconds. For this host, as can be seen in Fig. 7, the file (in this case 84 Mbytes) was transferred in 5-6 seconds.

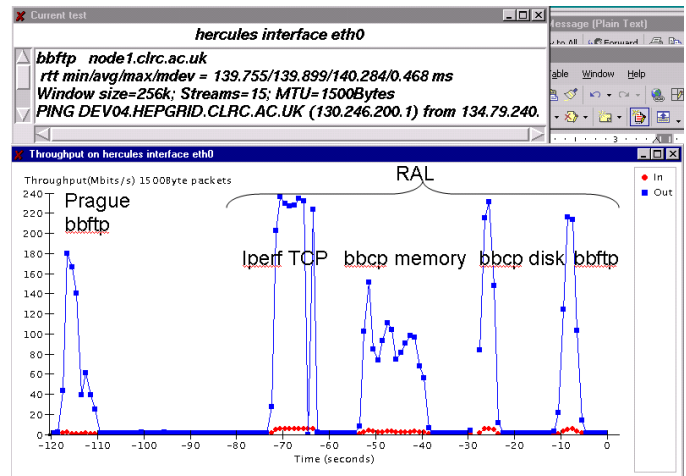


Figure 7: UTH/ifconfig plot of sequential throughputs.

E. Flood iperf Throughput

To see how much aggregated throughput we could achieve we modified the sequential tests to run iperf in TCP mode simultaneously and continuously to 4 groups of 5-7 remote hosts/group. We referred to this as the “flood” mode. We set up 4 monitoring hosts (stier, haan, hp3, and hp4; keeshond continued to run in sequential mode) in flood mode. Each monitoring host used NIC GE I, and sent the iperf TCP data to a different group of remote hosts. The sequential iperf/TCP throughputs to these remote hosts from keeshond at iGrid2002, via the routes shown in Fig 3, are shown in Fig 8. The members in each group were chosen to make the aggregate throughputs for each group about equal. The sustained throughputs achievable in flood mode estimated from the UTH/ifconfig plots on stier, are also shown in Fig. 8, together with which group was assigned to which monitoring host. We had two demonstration periods when we could send as much data as possible without regard for others.

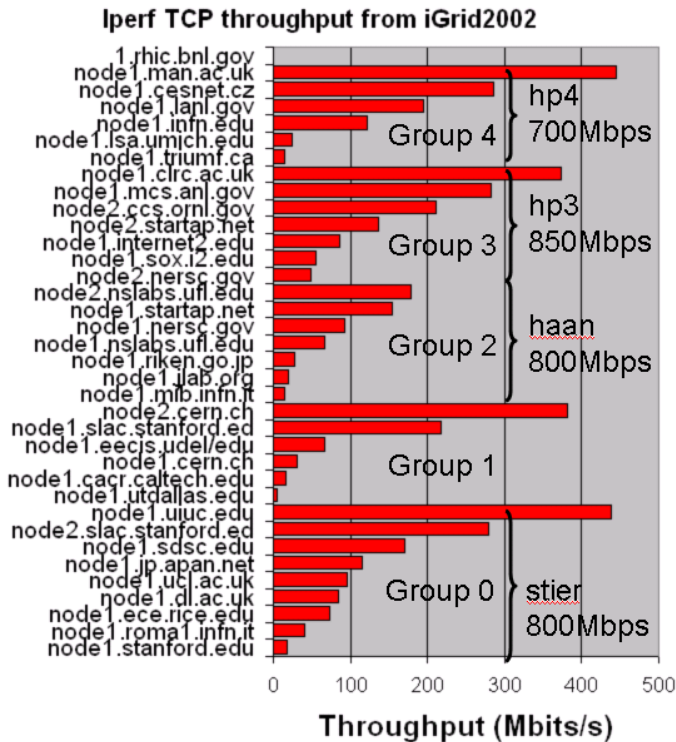


Figure 8: Iperf TCP sequential throughputs for remote host groups measured from keeshond at iGrid2002. The aggregate throughputs measured from stier are also shown, together with which group was assigned to which monitoring host.

The throughput observed in the iGrid2002 Cisco 6509 switch for the port attached to stier for the demonstration on Thursday 26 September can be seen in Fig. 9. The peaks in throughput correspond well with our demonstration time from 9am to 12 noon. These points were averaged over 5 minute intervals (UTH/ifconfig was averaged over 2 second intervals), and indicate that we achieved up to 675Mbits/s peaks and over 600 Mbits/s for sustained intervals. This is ~20% lower than we measured using UTH/ifconfig.

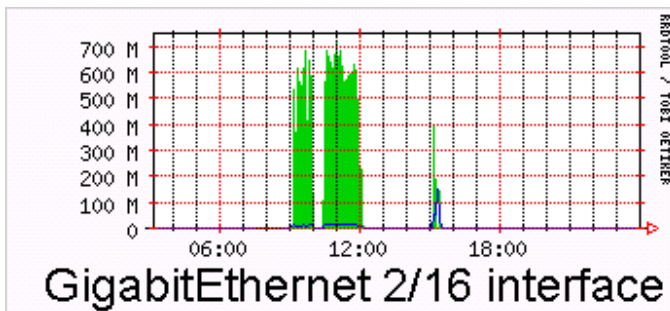


Figure 9: iGrid2002 switch throughput measured for the interface to stier on Thursday 26th September. The left hand axis is in bits/s, i.e. 0 - 750Mbits/s.

Due to lack of time, we were unable to make UTH/ifconfig run on hp3 and hp4. Thus when we ran the flood iperf demonstration to achieve the maximum throughput, we had to rely on throughput data from the iGrid2002 Cisco

6509 switch. Fig. 10 shows the Cisco 6509 throughput from 6am Wednesday 25th to 4pm Thursday 26th September 2002. We turned on the flood iperf throughput demonstration at around 9am Thursday. It is seen that during this time the outgoing traffic (the line representing traffic going from iGrid2002 to SURFNet) increases from a background of 1 Gbits/s to 3 Gbits/s.

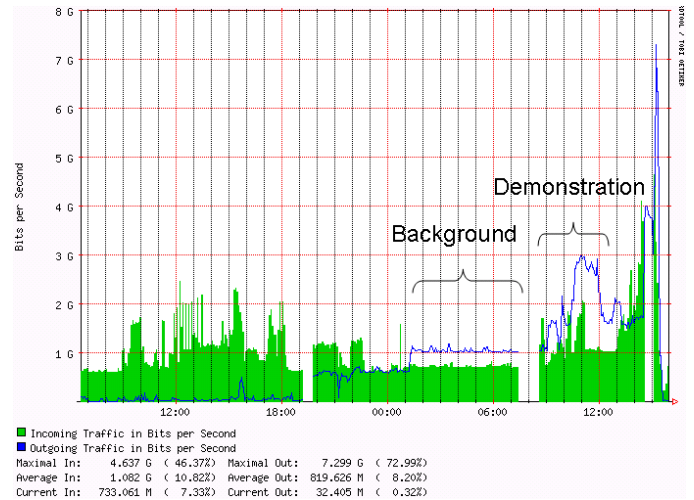


Figure 10: iGrid2002 router throughput to SurfNet.

IV. EXPERIENCES & RESULTS

A. PingWorld

The ping based PingWorld is low network impact (< 400 bits/s sent to each remote host), it requires no installation since the client is a Java applet (it may be necessary for some Web tools to be installed on the client), and the server comes pre-installed with most modern hosts), it provides RTT, loss and availability. It is particularly useful for monitoring losses to remote hosts with poor network performance (e.g. hosts in developing countries). It does not provide a good way to measure throughput for high performance paths since typically not enough pings are sent to measure < 1% losses, and even if they were, the ping losses are unlikely to accurately match TCP losses, since TCP induces congestion losses.

B. Iperf and Application Throughputs

We made measurements of TCP throughput and file copying to 34 hosts in 10 countries. The achievable iperf TCP throughputs (see Fig. 1) varied by more than a factor of 10 from site to site. By design, the hosts with 1000 NICs had higher speed connections (typically 622Mbits/s or higher) to the Internet and hence higher performance was observed.

Iperf TCP throughputs achievable to U.S. hosts, well connected (≥ 155 Mbits/s) to the Abilene and ESnet networks, were around 70-120 Mbits/s (see ORNL, LANL, NERSC, SLAC, SDSC, Florida, Internet2). ANL & UIUC had even higher achievable throughputs. ANL had a direct connection to StarLight. We are investigating why UIUC performed so well. There were 3 exceptions that had low performance: the SOX

host was unable to support parallel streams of TCP data, and so was only able to achieve 44 Mbits/s on average; the iperf TCP performance of the Michigan remote host seen from NIKHEF (and from SLAC) has recently doubled with no changes in streams or window size. Further investigation is in progress; the poor performance of the Caltech host from iGrid2002 compared to from SLAC (220 Mbits/s) is still under investigation. All these hosts measured more throughput from SLAC, so the throughput was not limited at the remote end. The other U.S. lower performance hosts connected to Abilene and ESnet all have only 100Mbit/s NICs and in some cases their site connectivity is limited to 43Mbits/s or less.

The APAN/Japan host achieved ~100Mbits/s on average from iGrid2002, whereas seen from SLAC it achieves 200-400Mbits/s. This is similar to the throughput achieved to most well connected US hosts on Abilene and ESnet.

European well connected hosts (see Manchester, RAL, Prague) achieved throughputs of 250-450 Mbits/s on average, or 2-3 times that achieved by similar US and Japanese hosts. These European hosts also see similar throughput from SLAC. The asymmetry in the achievable throughputs from NIKHEF to US hosts and SLAC to European hosts is worthy of further investigation.. The poor performance to IN2P3 was probably caused by unusual routing during iGrid2002 and lack of time to optimize the windows. After adjusting the windows and streams we are now achieving about 200 Mbits/s from NIKHEF (the site for iGrid2002) to IN2P3.

Performance of disk-to-disk file copies (bbcp disk-to-disk and bbftp) for high performance links (> 100Mbits/s) tended to be well below that of iperf TCP. From other work we believe this is due to disk and file performance issues. Typically the disk-to-disk file copies max out at 50-100Mbits/s. Comparing throughputs reported by bbcp disk-to-disk vs. bbftp indicates that in most cases bbftp is slower. However, this is an artifact due to bbcp starting the timer when the sockets are all set up, whereas bbftp starts the timer when the program is entered. The former is more accurate for measuring network throughput, the latter for measuring what the user may expect. We are modifying IEPM-BW and bbcp to report the throughput by both methods.

C. Bandwidth Estimation

The ABWE tool is able to quickly show short term (within a couple of minutes) changes in network performance (e.g. due to route changes or big changes in traffic and congestion) or a host being unreachable (shows as a flat line) in real-time. During iGrid2002, to make the graphs show trends more clearly, we displayed Exponentially Weighted Moving Averages (EWMA) of the measurements, i.e. the current average avg_i is given by:

$$avg_i = (1 - a) * y_i + a * avg_{i-1} \quad (1)$$

where y_i is the current measurement, and a is a constant that was set to 0.9 during the demonstration. This provides fairly heavy smoothing, and since the plots were updated once a minute, this meant that changes would not be visible for a few minutes (depending on the magnitude of the change). We also

have a version of ABWE that presents both the actual value and the EWMA. This plot is noisier, but allows one to see changes more quickly. Fig 11 shows an example of a sudden change in network performance for CERN without and with EWMA smoothing with $a = 0.9$. It can be seen that the sudden changes in throughput are seen within 2 minutes (each point is for one minute) for the unsmoothed data and after a few minutes for the smoothed data. Further investigation showed that the first drop was a change from the normal route from SLAC to CERN going via StarLight to a new route shared with commercial traffic going via BBNPlanet. The increase, at 280 minutes, was caused by the route to CERN changing to going via GEANT.

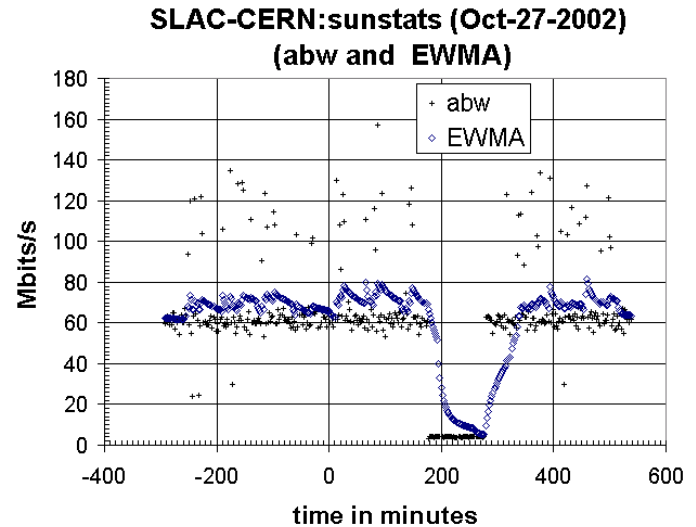


Figure 11: Unsmoothed ABWE (abw) measurements vs EWMA smoothed values.

We compared the SLAC iperf TCP throughputs averaged over 60 days from August 24 through October 26 2002 with the SLAC ABWE results. The correlation was strong (square of the correlation coefficient, $R^2 > 0.6$), see for example Fig. 12. In Fig. 12 each point represents a single host with the x value being the 60 day average of the iperf TCP measurements, and the y value being the average ABWE for a four hour period on October 26, 2002 (the pattern of the ABWE values stay fairly consistent from weekday to weekday, so the actual choice of time period has little effect). The line is a fit to a straight line constrained to go through the origin. We also noted that if there were large (e.g. diurnal) variations in the iperf TCP measurements they also showed up in the ABWE measurements. This is encouraging, and we are studying it further. If it bears out, then we hope to be able to use the more heavyweight iperf TCP measurements to normalize the more frequent ABWE measurements. Besides providing low network impact (as configured at iGrid2002, each bandwidth measurement took about 60 kbits/remote host, and the measurements were repeated at 1 minute intervals, so the average bandwidth/remote host was 1 kbits/s) real-time presentations of network performance, we have also found the ABWE measurements valuable for identifying hosts that may need their window and stream settings to be optimized.

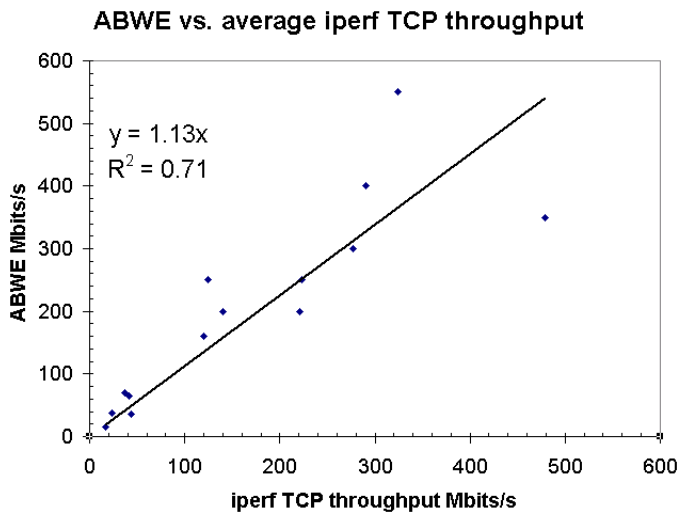


Figure 12: Correlation between ABWE measurements and the average iperf TCP throughputs.

D. Flood iperf Throughputs

We had two demonstration time-slots when we were authorized to send as much data as we wanted. During both of these periods we were able to send just over 2 Mbits/s from our monitoring hosts to up to 34 remote hosts in 10 countries.

V. CONCLUSIONS

We demonstrated a wide range of active Internet end-to-end measurement tools that cover a wide range of needs from monitoring low performance paths to high performance paths, and from long term to real-time.

The mix of a low network intrusive packet pair bandwidth measurement tool (ABWE) with a more intrusive, user-centric TCP throughput tool (iperf) appears promising to provide low-impact short term (updates/minute) real-time measurements with good normalization.

We have observed that with standard configurations (operating system, TCP stack, MTU, best effort traffic only with no special quality of service) over well-provisioned networks, we can achieve over 400 Mbits/s host-to-host over trans-Atlantic distances. We were also able to achieve 800-900 Mbits/s from a single host to a group of 6 remote hosts. To do this requires a careful choice of common off-the-shelf hardware (GHz cpu, high speed buses, e.g. 66MHz 64bit PCI, GE NIC), large TCP buffer/window sizes and multiple parallel TCP streams. We look forward to testing with new TCP stacks (see for example [33]) to evaluate their performance on an extensive set of paths, and to see whether we can remove the need for multiple streams.

ACKNOWLEDGMENT

We gratefully acknowledge the assistance of Jerrod Williams with the IEPM-BW toolkit, Paula Grosso for the UTH/`ifconfig` plots and Tony Johnson with the PingWorld application, the iGrid2002 organizers, and NSF for funding the travel. We thank Cees de Laat and NIKHEF for

providing the measurement hosts at iGrid2002, and all the administrators and contacts of the remote hosts for providing the hosts and responding to questions and requests.

REFERENCES

- [1] "BaBar Collaboration Home Page". Available <http://www.slac.stanford.edu/BFROOT/>
- [2] "Jefferson Lab Home Page". Available <http://www.jlab.org/>
- [3] "The Collider Detector at Fermilab". Available <http://www-cdf.fnal.gov/>
- [4] "The D0 Experiment". Available <http://www-d0.fnal.gov/>
- [5] LHC – Large Hadron Collider. Available at <http://lhc-new-homepage.web.cern.ch/lhc-new-homepage/>
- [6] "The DataGrid Project". Available <http://eu-datagrid.web.cern.ch/eu-datagrid/>
- [7] "Particle Physics Data Grid". Available <http://www.ppdg.net/>
- [8] "GriPhyN – Grid Physics Network". Available <http://www.griphyn.org/index.php>
- [9] "iGrid2002". Available at <http://www.igrid2002.org/>
- [10] "Atlas hierarchical Computing Mode". Available <http://www-iepm.slac.stanford.edu/monitoring/bulk/igrid2002/tier.gif>
- [11] "Bandwidth Challenge from the Low-Lands". Available <http://www-iepm.slac.stanford.edu/monitoring/bulk/igrid2002/>
- [12] "iGrid2002 Featured Networks". Available <http://www.igrid2002.org/>
- [13] "Bandwidth to the World. – Host Requirements". Available <http://www-iepm.slac.stanford.edu/monitoring/bulk/sc2001/host-req.html>
- [14] R. L. Cottrell. "Network Monitoring Tools". Available at <http://www-iepm.slac.stanford.edu/monitoring/bulk/sc2001/host-req.html>
- [15] "CAIDA Tools Taxonomy". Available <http://www.caida.org/tools/>
- [16] "Ping". Available <http://www.webopedia.com/TERM/P/PING.html>
- [17] "The PingER Project". Available <http://www-iepm.slac.stanford.edu/pinger/>
- [18] "Introduction to the NLANR AMP Project for HPC sites". Available <http://watt.nlanr.net/AMP/>
- [19] "Introduction to the Surveyor Project". Available <http://www.advanced.org/surveyor/>
- [20] "RIPE traffic measurement project". Available <http://www.ripe.net/test-traffic/Documents/RIPE/RIPE179/note.html>
- [21] V. Jacobson, "traceroute". Available <ftp://ftp.ee.lbl.gov/traceroute.tar.z>, 1989.
- [22] M. Jain, C. Dovrolis, "End-to-End Available Bandwidth Methodology, Dynamics, and Relation with TCP Throughput". ACM SIGCOMM 2002, Pittsburgh PA, August 2002
- [23] Network Characterization Service (NCS). Available <http://www-didc.lbl.gov/pipechar/>
- [24] A. Tirumala, F. Quin, J. Dugan, J. Ferguson and K. Gibbs, "Iperf 1.6 – The TCP/IP bandwidth Measurement Tool". Available <http://dast.nlanr.net/Projects/Iperf/>
- [25] A. Hanushevsky, A. Trunov, R. L. Cottrell, "P2P Data Copy Program bbcp". Available <http://www.ihep.ac.cn/~chep01/presentation/7-018.pdf>
- [26] "Large files transfer protocol". Available <http://doc.in2p3.fr/bbftp/>
- [27] The GridFTP Protocol and Software". Available <http://www.globus.org/datagrid/gridftp.html>
- [28] M. Mathis, J. Semske, J. Mahdavi, and T. Ott. "The macroscopic behavior of the TCP congestion avoidance algorithm". *Computer Communication Review*, 27(3), July 1997.
- [29] "PingWorld". Available <http://jas.freehep.org/demos/PingWorld/>
- [30] "UTH". Available at <http://jas.freehep.org/documentation/UTH/>
- [31] "SLAC WAN Bandwidth Measurement Tests for Monitoring Site iGrid2002, NL". Available at http://www.dutchgrid.nl/datatag/html/slac_wan_bw_tests.html
- [32] R. L. Cottrell and C. Logg. "Experiences and Results from a New High Performance Internet Monitoring Infrastructure", SLAC-PUB 9202.
- [33] Sally Floyd, "HighSpeed TCP for Large Congestion Windows", draft-floyd-tcp-highspeed-00.txt, July 200.

