

Analyzing Traffic across the Greek School Network

Costas Kattirtzis, Emmanuel Varvarigos, Kyriakos Vlachos, *member IEEE*,
George Stathakopoulos and Michael Paraskevas

Abstract— In this paper, we present a comprehensive traffic analysis of the Greek School Network (GSN), a wide area network designed to provide Internet access and services to about 15,000 units of primary and secondary education schools and administration offices. In our analysis, we have used measurements from the PATRAS region node obtained through the Cisco NetFlow and FlowScan tools. We have used classical analysis to obtain protocol and application traffic statistics. Our study revealed that TCP traffic is dominant in the network, while nearly 50% of the outgoing and 37% of the incoming traffic is Peer-to-Peer (P2P) traffic, with a further 25.6% of traffic using not registered ports and suspected to be P2P as well. Finally, we have also observed a remarkable traffic locality phenomenon in the P2P services, where more than the 90% of the traffic was heading or generated by 50 hosts.

Index Terms— Peer-to-Peer traffic, traffic locality, traffic measurements

I. INTRODUCTION

Internet evolution has been accompanied by the growth of traffic volume, the development of new protocols, and the proliferation of a variety of new Internet access technologies. Network traffic measurements are useful for network troubleshooting, workload characterization, and network performance evaluation. In order to detect the invariants in a rather dynamic traffic structure, measurements and analysis of genuine network traffic traces are important and continuing tasks. However, the traffic carried by Internet backbones presents unpredictably variable and complex patterns [1]. The complexity of the Internet is mainly caused by the aggregation of traffic flows from many end systems, users and applications [2]. Previous studies of TCP/IP traffic have examined the statistics of the aggregated packet arrival processes in local-area [3] and wide-area networks [4]-[6]. These studies have shown that packet inter-arrival times do not produce a smooth “Poisson like” superposition process but rather follow a packet train model. Other studies present a number of analytic models for describing the traffic characteristics of telnet, NNTP, SMTP and FTP connections [7], and others, show that

This work was supported in part by the Ministry of Education and Religious Affairs of Greece.

The authors are with Research Academic Computer Technology Institute, 61 Riga Ferraïou str., Patra, Greece.

C. Kattirtzis, E. Varvarigos and K. Vlachos are also with Computer Engineering and Informatics Department, University of Patras, GR26500, Rio, Patra, Greece (tel: +30 2610 996990, fax: +30 2610 960350, email: kvlachos@ceid.upatras.gr).

wide area traffic processes described by the Poisson model are valid only for modeling the arrival of user sessions and that Ethernet traffic is statistically self-similar [5].

In this paper, we present a study of traffic patterns and flow profiling statistics of a metropolitan area network, of the Greek School Network -GSN (www.sch.gr) in the PATRAS prefecture. The GSN is a wide area network that connects all primary and secondary education schools in Greece. Our study revealed that TCP protocol is dominant in the network and accounts for 95% of the bytes, 61.7% of the flows and 84.1% of the packets, while the UDP protocol accounts for only 4.4%, 14.5% and 34.5% of the bytes, flows and packets respectively. Furthermore, we have found out that 87% of the flows carry only 5-12 packets and less than 1% of the flows contain more than 12 packets. A remarkable observation is that nearly 50% of the outgoing and 37% of the incoming traffic is Peer-to-Peer (P2P) traffic with a further 25.6% of traffic using unregistered ports and which is suspected to be P2P traffic as well. At 5-min scales, the traffic displays a distinctive piecewise-linear non-stationarity, together with evidence of long-range dependence. This seems to be in general agreement with recent theoretical models for large-scale traffic aggregation. Finally, our analysis revealed a strong traffic locality phenomenon (most probably caused by P2P services), where more than the 90% of the traffic was heading or generated by 50 hosts.

The remainder of the paper is organised as follows. Section II presents an overview of the GSN and the measurement methodology, Section III presents basic metrics and traffic statistics of the MAN while Section IV concludes the paper.

II. NETWORK ARCHITECTURE AND MEASUREMENT METHODOLOGY

The GSN is the educational intranet of the Ministry of Education and Religious Affairs (www.ypepth.gr). The distribution network of GSN has been operational for 5 years now, providing network connectivity and Internet services to all the schools and administration offices across all the fifty one (51) prefectures of Greece. The GSN is hierarchically structured into three layers in order to manage the complexity that comes with the large number of sites that are covered. The Greek Research and Technology Network (GRNET, www.grnet.gr) is used as its core network, with which it interconnects at seven main points. The distribution network consists of network and computational equipment installed at the capital of every prefecture and in that way ensures optimal

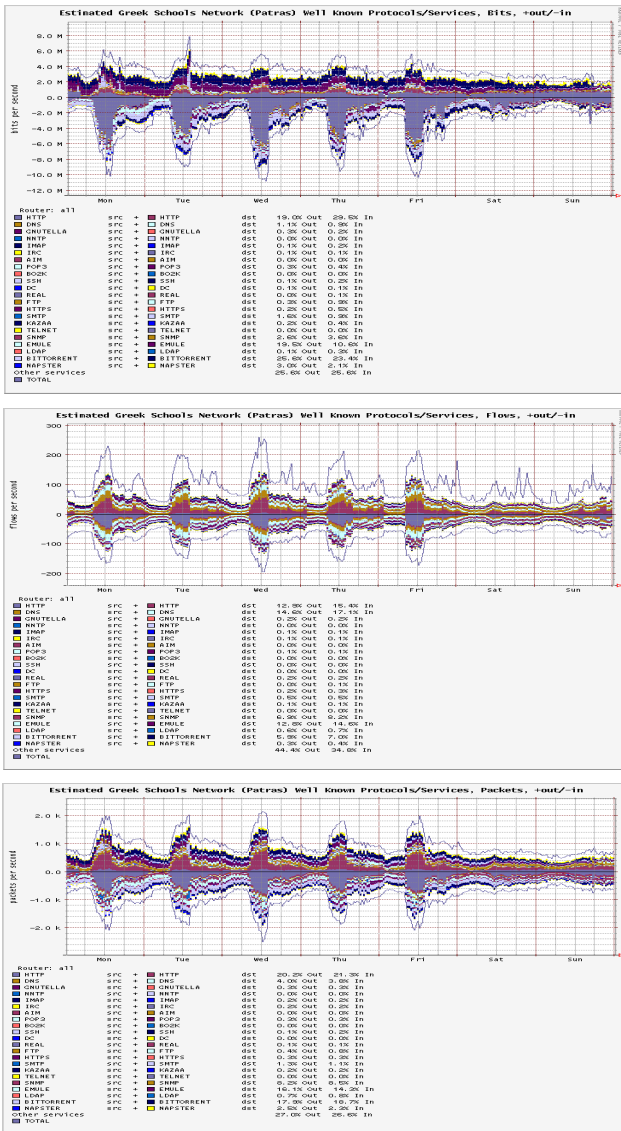


Figure 2: Traffic load distribution of (a) bits, (b) flows and (c) packets by services for ingoing and outgoing traffic

Figure 3 demonstrates the expected predominance of small-sized packets in the traffic. Almost half of the packets have length that is less than 64 bytes, while 66.2% of them are shorter than 352 bytes and 70% are shorter than the typical TCP MSS of 576 bytes. It must be noted here that packet size distribution in the GSN is not tri-modal as described in [4] but rather dual-modal as we have observed from the cumulative density function of the packet size distribution. This is because the amount of 552-byte or 576-byte MTU packets in the network is not large.

The large peak exhibited at small packet sizes is primarily caused by TCP control segments, which are 40 bytes long (SYN, FIN, RST packets), and by HTTP protocol having 41 bytes length and carrying single characters. Concerning the first sample, the 16.76% of the traffic is in the 65-192 bytes range and is mainly caused by P2P applications, while 2.1% of the traffic is in the range of 545-576 bytes and is due to the TCP mechanism for minimizing packet fragmentation. In the

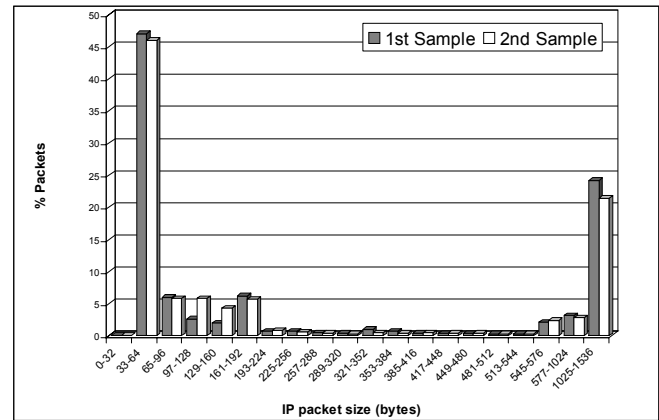


Figure 3: Packet size histogram for two samples monitored on our link

GSN network, Path MTU Discovery is employed with MTU size of 1500 bytes and thus few packets exist in the range of 552 to 576 bytes. The 3.14% of the traffic is in the 577-1024 bytes range and is mainly caused by P2P traffic. Large packets of 1025-1536 bytes length correspond to 24.18% of the total traffic and are comprised of full Ethernet frame packets, 1500 bytes long, stemming mainly from P2P services like eMule and BitTorrent [12]. It must be noted here that the percentages presented above have a rounding error of $\pm 0,5\%$, adding to 1% error after summation of the measurements. This rounding error is negligible and does not affect the measurements of our study.

Figure 4 and Figure 5 show box-plot graphs of the distribution of the packet size and the number of packets per flow for several well known internet applications. In Figure 4 the Y-axis shows the number of packets per flow for various applications. The X-axis corresponds to the mean number of flows observed for the first sample. In the boxplot graph, the outside ends indicate the maximum/minimum of the mean weekly value, while the box itself shows the middle half of the data and the diamond in the middle the median of the number of packets per flow for each application. As expected, the mean weekly number of flows varies significantly from week to week for applications that produce a small number of flows per week. From these two figures, useful information regarding the type of the transfer can be derived. For example, data transfers can be categorized in three distinct groups: interactive, transaction oriented and bulk transfers, as it is proposed by K. Claffy et al. in [13]. A quarter of the applications in the monitored link is interactive and only sends end-to-end packets carrying payloads that consist of a single or only a few characters. Applications such as ICMP, IGMP, TCP-Telnet, UDP-DNS and UDP-NTP belong to this group. Applications like TCP-FTP, TCP-SMTP, UDP-TFTP and IPINIP are transaction-type applications and include a relatively large number of bytes per packet. In contrast to these two categories, other applications like TCP-X, TCP-FTPD, TCP-Frag, TCP-WWW, GRE belong to the bulk data transfer-style and usually send very large packets. These applications do not necessarily dominate the traffic simply

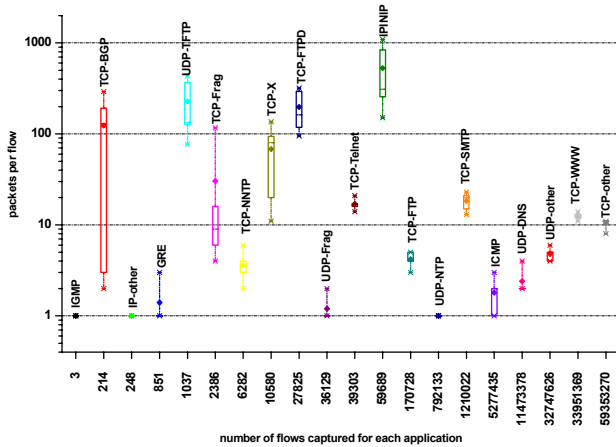


Figure 4: Distribution of mean weekly packets per flow by protocol type

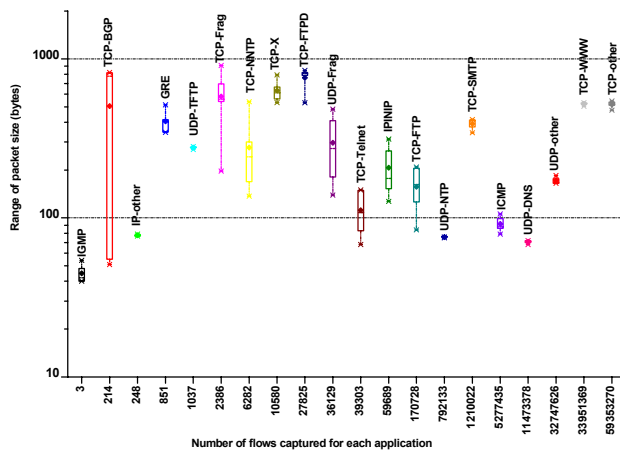


Figure 5: Mean weekly packet size per flow by protocol type

because applications like TCP-FTPD are rarely used and thus the corresponding number of flows is limited. On the other hand, interactive applications such as ICMP, UDP-NTP, UDP-DNS even though they exchange only a few packets (1-5 packets per flow), they dominate the traffic mix because they are widely used resulting in a large number of active flows. Combining these two graphs, we can conclude that the TCP applications, TCP-WWW and TCP-other and the UDP-other are the dominant applications with respect to the number of flows and the number of bytes transmitted.

B. Flow Analysis

The majority of the flows contain a relatively small amount of packets. In particular, with results stemming from Figure 4 and Figure 5, we have calculated the cumulative distribution function for the flow size. Almost 87% of the flows carry 5-12 packets and less than 1% of the flows contain more than 12 packets. Generally, our results show that the majority of the flows consist of only a few packets and there is a negligible amount of large flows. Furthermore, we have found that the flows produced by the FTP application last much longer than the other flows. TCP-WWW and TCP-other produced the majority of the flows. These flows had a mean duration of 6.14 and 6.36 seconds respectively and from this we can

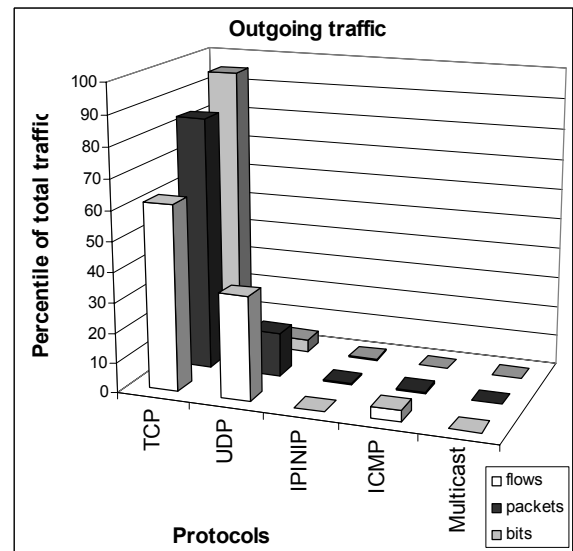


Figure 6: Traffic composition by protocol for the outgoing traffic for a week sample

conclude that the majority of the flows have a small duration.

C. Protocol Analysis

Figure 6 reveals the composition of the traffic of the second sample. We observe that TCP traffic by far dominates the traffic mix. Over the period of a week, TCP averages 95% of the bytes, 61.7% of the flows and 84.1% of the packets on the monitored outgoing link. UDP is the second largest category at about 4.4% of the outgoing traffic and 5.2% of the incoming traffic in terms of bytes. The other IP protocols plotted, Multicast, IPINIP and ICMP individually make up a negligible percentage of the overall traffic and for this reason we are only going to discuss TCP and UDP protocols. It must be noted here that multicast traffic is not supported by the GSN network and thus there is no outbound multicast traffic. Furthermore, comparing the graph of bits per second with the graph of flows per second and packets per second, we can detect that the size of the incoming packets is much larger than the size of the outgoing packets. This is mainly caused by the packets of the HTTP traffic. It is useful to see and compare TCP and UDP traffic. UDP constitutes 4.4% of the bytes, 14.5% of the flows and 34.5% of the packets. Using the evidence from the above graph we can deduce that TCP flows have very different characteristics than the UDP ones. More specifically, TCP uses more and larger packets per flow than UDP. Comparing our results with the ones in [14] we can conclude that the composition of the traffic observed in a School network is more or less the same with the traffic observed to the real world networks.

D. Service Analysis

In this section, we analyze traffic in terms of well known applications over a week period and we also discuss the relationship between these applications and the underlying protocols. As previously noted, Figure 2(a) to (c) displays the traffic load distribution of bits, flows and packets per second by service type for the ingoing and outgoing traffic. It can be

deduced that, in terms of bytes transferred, almost 50% of the outgoing traffic and nearly 37% of the incoming is P2P traffic, while 19% of the outgoing traffic and almost 30% of the incoming traffic is HTTP traffic. There is a high percentage of “unknown” services, close to 25.6% of the outgoing and the incoming traffic. These types of services are spread among a range of TCP and UDP port numbers that are not registered to IANA.

SNMP is in the fourth position in terms of contribution to the traffic load, representing approximately 2.6% and 3.6% of the incoming and outgoing traffic, respectively. SMTP follows with 1.6%, DNS with 1.1% and POP3 with 0.3% of the outgoing traffic. All the remaining applications, such as NNTP, IMAP, IRC, BO2K, SSH, Real, HTTPS, Telnet, LDAP, AIM, represent only 0.6% and 1.4% of the outgoing and the incoming traffic, respectively.

Based on the same results we can derive useful information regarding the daily distribution of the traffic volume. For example, regarding the HTTP (web) application, we observed that the profile of its daily load distribution fits closely the corresponding profile of the TCP protocol. This was expected, since HTTP services mainly make use of this protocol. It is worth noting here, that the outgoing packets of HTTP services are about two times smaller than the incoming packets, even though the fraction of packets and flows is almost the same on the incoming and outgoing link. This is because the outgoing packets contain TCP control datagrams and request packets, which are short in length. The daily mean of flows per second (fps) over the week varies from 36 to 260 for the outgoing traffic and from 38 to 195 for the incoming traffic. Similarly the mean bit rate varies from 1.75 to 7.9Mbps and from 1 to 10.8Mbps and the mean packet rate from 510 to 2200 and from 430 to 2560 for the outgoing and incoming traffic, respectively.

Regarding P2P services, they generate 24hours per day a prominent amount of traffic that constitutes up to 50% of the total traffic. During the last few years these applications have emerged as the predominant bandwidth-consuming services of the Internet. Capturing and analyzing P2P services is not an easy procedure. In our measurements, relying on previous knowledge of P2P networks, we extend the configuration of CUFlow to capture a larger percentile of P2P traffic. Although a large fraction of P2P traffic was not recognized, and thus P2P traffic was underestimated, the percentile of the captured P2P traffic is high enough to derive useful information and analyze its effect on the School network.

During the last 7 months P2P services contributed to the total outgoing traffic at a percentage ranging from 32.3% to 48.7%, and to the total incoming traffic at a percentage ranging from 14% to 30.9% verifying the findings of Karagiannis et al. [12], [15]. Regarding the traffic pattern produced by the P2P services during the day, it can be observed from Figure 2 that P2P applications are 24 hours active and are only reduced during the night, but not as much as the other services. A significant characteristic of P2P applications is that they do not follow the traffic pattern of the

remaining aggregated data traffic, and they do not surge during the peak hours.

eMule and BitTorrent were by far the two most prevalent protocols, while Gnutella, Kazaa and Direct Connect made up individually an inconsiderable percentage of the overall traffic. Comparing the findings for BitTorrent and eMule, we can deduce that BitTorrent flows are larger in size and they contain less but larger packets (in bytes) than eMule flows, while BitTorrent produces the biggest flows of all the services.

It worth noticing here that newer samples captured after the 19th of December 2004 showed a very steep reduction in the use of BitTorrent. More precisely, in a sample captured in March 2005, the use of BitTorrent service was dramatically reduced, and was responsible only for the 3.5% and 2.4% of the outgoing and incoming traffic, respectively. This was not a coincidence. On that day Suprnova.org –one of the biggest sites offering “*torrent*” meta files– was shut down by the owners due to legal warnings. This had the additional side effect that the unrecognized traffic in the school network was increased and thus we can assume that BitTorrent users are now employing other untraceable P2P applications.

Concerning the arrival rate of P2P applications, striking results were observed. By aggregating the data in 5-minute time bins we have observed that the slopes of the cumulative traffic curves of bytes, flows and packets caused by these applications during a day interval remain relatively constant throughout the day. We have noticed the same pattern on a weekly interval with the distinction that we had an almost constant rate during the weekdays and a different rate (but relatively constant again) during the weekends. Similar measurements at multi-second time scales [16] have shown that the arrival rate may remain constant for several minutes at a time but it is clearly a time varying function with non-stationarity characteristics.

Measurements on the number of packets per flow revealed that the majority of P2P flows contain a relatively small number of packets. The biggest P2P flow captured in a day contained 21 packets, while 83% of the flows carried at most 12 packets. The average size of a flow was found to be 9 packets. As far as the type of the P2P flows, P2P applications belong to the bulk data transfer-style applications. 50% of the flows caused by these services are less than 5.5 Mbytes while the biggest one is 17.5 Mbytes. The mean P2P flow size was found to be 6.1 Mbytes, which is much bigger than the mean flow size of web traffic and other bulk data transfer services.

Mail protocols, such as SMTP and POP3, despite their practical importance, produce only a small amount of traffic, namely, 0.6% of the flows, 1.6% of the packets, and 1.9% of the bytes. Their use increases during the peak hours, albeit they do not seem to exhibit a specific daily traffic pattern.

SNMP and DNS use UDP as the underlying protocol and thus both services, produce a large fraction of the flows, a small fraction of the packets and an even smaller fraction of the bytes transferred. DNS accounts for 14.6% of flows (8 to 65 fps), 4% of packets (21-190 pps) and only 1.1% of bytes

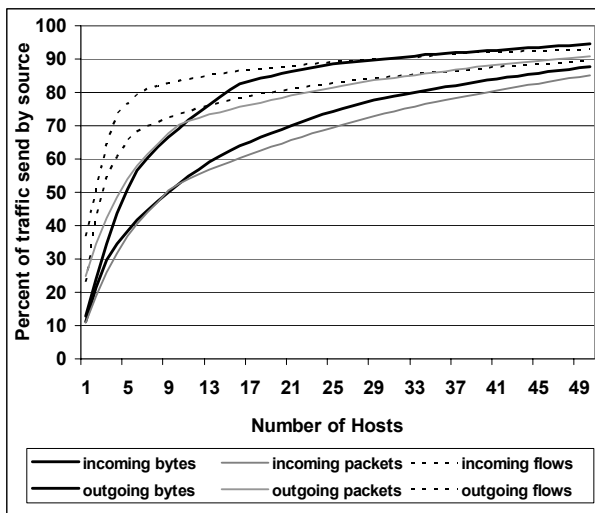


Figure 7: Cumulative distribution function of traffic to favorite sources and destinations

(22-200 bps). The daily traffic size pattern (volume of bytes) of DNS, exhibits the same pattern as the aggregated data. Concerning SNMP, it exhibits the same pattern with the incoming/outgoing traffic. This was expected since SNMP protocol sends control messages from the network devices to a central workstation console (NMS) and vice versa to control and monitor network resources. SNMP service in GSN has been programmed to operate during the early hours so as not to consume network resources during the peak hours. The max traffic rate was captured in the morning, with 12fps, 174pps and 240kbps.

Finally, with respect to FTP services, no specific traffic pattern was detected. As expected, FTP has a higher percentile of bytes and packets than flows, meaning that FTP flows last long, each carrying a large size of traffic.

E. Traffic Locality

Traffic locality is a special case of the locality phenomenon seen in computer systems. The traffic load in terms of flows, packets or bytes is not randomly distributed to all source and destination addresses. The aggregated traffic is frequently dominated by a small number of transmitting hosts, both in short and long time scales. Figure 7 shows the 50 most busy sources and destinations in a 5-minute sample, and it indicates that these 50 hosts (out of the 6188 in the PATRAS prefecture) are responsible for the generation of almost 95% of the bytes, 93.1% of the flows and 90.9% of the packets. The same users consume also a large fraction of the incoming traffic namely 76.6%, 77.5% and 52.5% of the incoming bytes, packets and flows, respectively.

A similar observation was also made for the 50 hosts that consume the majority of the incoming traffic. In particular, these hosts consume 87.8% of the bytes, 89.5% of the flows and 85.2% of the incoming packets. Specific site pairs show even more remarkable locality. The first 5 busiest destinations (incoming traffic), absorb 38.2% of the bytes, 33.1% of the packets and 34.2% of the flows. Similarly, the first 5 busiest

sources produce 51% of the bytes, 20% of the packets and 24.2% of the flows. It is worth noticing here that the locality phenomenon for the specific sources/destinations hosts remained unchanged even in longer time scales of 250 minutes.

IV. CONCLUSION

In this paper, a comprehensive traffic analysis of the metropolitan area network of the Greek School Network in the PATRAS prefecture was presented. Our analysis revealed that TCP traffic dominates in the network, while nearly 50% of the outgoing and 37% of the incoming traffic is P2P traffic with a further 25.6% of traffic using unregistered ports and which is suspected to be P2P traffic as well. At 5-min scales, P2P traffic displays a distinctive piecewise-linear non-stationarity, while the mean packet and byte size of a flow was found to be 9 packets and 6.1MB respectively. Finally, our analysis revealed that strong traffic locality exists, most probably caused by P2P services, where more than 90% of the traffic was heading or generated by only 50 hosts. It is evident that P2P accounts for a considerable fraction of the aggregated traffic, and this has to be taken into account for future network extensions.

REFERENCES

- [1] W. Leland, M. Taqqu, W. Willinger, D. Wilson, "On the self-similar nature of Ethernet traffic", IEEE/ACM Transactions on Networking, Feb.1994.
- [2] V. Paxson, "Growth Trends in Wide-Area TCP Connections", IEEE Network, July 1994.
- [3] R. Jain and S. Routhier, "Packet Trains — Measurements and a New Model for Computer Network Traffic," IEEE JSAC, 4(6), pp. 986-995, Sep. 1986.
- [4] K. Thompson, G. Miller, R. Wilder, "Wide-Area Internet Traffic Patterns and Characteristics", IEEE Network Magazine, Nov. 1997.
- [5] V. Paxson, S. Floyd. "Wide-Area Traffic: The Failure of Poisson Modeling" IEEE/ACM Transactions on Networking, June 1995.
- [6] K. C. Claffy, G. Miller, and K. Thompson, "The nature of the beast: Recent traffic measurements from an Internet backbone," in Proc. INET'98, Geneva, Switzerland, July 1998
- [7] V. Paxson, "Empirically-Derived Analytic Models of Wide-Area TCP Connections," IEEE/ACM Transactions on Networking, 2(4), Aug. 1994.
- [8] NetFlow Services and Applications, Cisco White Paper. Available: http://www.cisco.com/warp/public/cc/cisco/mkt/ios/netflow/tech/napps_wp.htm
- [9] D. Plonka, "Flowscan: A Network Traffic Flow Reporting and Visualization Tool," Proc. 2000 USENIX LISA.
- [10] Daniel W. McRobb, "cflowd configuration", 1998-1999
- [11] Dave Plonka, "RRGrapher - the Round Rober Grapher, a Graph Construction Set for RRDtool"
- [12] T. Karagiannis, A. Broido, N. Brownlee, K. Claffy, and M. Faloutsos. "File-sharing in the Internet: A characterization of P2P traffic in the backbone". Technical report, 2004.
- [13] K. Claffy, G. Polyzos, and H.W. Braun, "Traffic Characteristics of the T1 NSFNET Backbone", Proceedings of INFOCOM '93, San Francisco, March, 1993.
- [14] M. Fomenkov, K. Keys, D. Moore, K. Claffy, "A longitudinal study of internet traffic from 1998-2003", WISICT, Cancun, January, 2004.
- [15] T. Karagiannis, A. Broido, N. Brownlee, K. Claffy and M. Faloutsos. "Is P2P dying or just hiding?" In *IEEE Globecom 2004 - Global Internet and Next Generation Networks*, 2004.
- [16] T. Karagiannis, M. Molle, M. Faloutsos, and A. Broido "A Nonstationary Poisson View of Internet Traffic," In IEEE INFOCOM '04, Computer and Communications Societies Conference on Computer Communications, Hong Kong, March 2004