

Fast Classification and Estimation of Internet Traffic Flows

Sumantra R. Kundu, Sourav Pal, Kalyan Basu, and Sajal K. Das

Center for Research in Wireless, Mobility and Networking (CReWMaN)
The University of Texas at Arlington, TX 76019-0015
{kundu, spal, basu, das}@cse.uta.edu

Abstract. This paper makes two contributions: (i) it presents a scheme for classifying and identifying Internet traffic flows which carry a large number of packets (or bytes) and are persistent in nature (also known as the *elephants*), from flows which carry a small number of packets (or bytes) and die out fast (commonly referred to as the *mice*), and (ii) illustrates how non-parametric Parzen window technique can be used to construct the probability density function (pdf) of the elephants present in the original traffic stream. We validate our approach using a 15-minute trace containing around 23 million packets from NLNR.

1 Introduction

There are two main aspects to the problem of Internet traffic flow characterization: (i) *how* to efficiently collect the flows, and (ii) how to accurately *infer* overall traffic behavior from the collected data. Due to limitations in hardware capabilities, it has been illustrated in [7] [22] how exhaustively collecting all packets of a flow does not scale well at high link speeds (OC-48+). Thus, current approaches to flow characterization are either based on: (i) *statistical sampling* of the packets [5][6], or (ii) *inferring* traffic characteristics primarily based on flows which carry a large number of packets (or bytes) and are long-lived in nature (i.e., the *elephants*) while ignoring flows which carry very small number of packets (or bytes) and are short-lived in nature (i.e., the *mice*) [10], or (iii) using appropriate *estimation algorithms* on lossy data structures (e.g., bloom filters, hash tables) [7][17] for recovering lost information. However, even in sampled traffic, separation of elephants and mice is a cumbersome task [8] since there exists no standard approaches to drawing the line between the two.

In this paper, we show that it is indeed possible to provide an analytical framework for identifying and classifying packets as elephants or mice by applying Asymptotic Equipartition Property (AEP) from Information Theory [18]. It is based on the observation that *all mice die young and are large in number*; while the proportion of elephants *is small in number* (around 1% – 2% of the traffic volume) and they have average longevity varying from a few minutes to days [1]. If the state space of the Internet flows is visualized to be an ergodic random process, then the existence of *typical sequence*, as defined by AEP, identifies the presence of elephants in the traffic volume. Such an approach requires

no prior knowledge of flow distribution, does not suffer from the side effects of false positives associated with Bayesian analysis, and involves minimal packet processing. We compare our approach with the well-known method of identifying packets as elephants based on the average flow longevity of greater than 15 minutes [8]. Our results from initial analysis on a single 15-minute traffic trace from NLNR [15] indicates that there exists a possibility that using definite values of longevity as cutoff limits for classifying flows as elephants might overestimate the frequency of occurrence of such flows. In the second part of the paper, we use a statistical non-parametric estimation technique based on the Gaussian kernel function for accurately estimating the density function of the underlying traffic, considering the probability density function (pdf) of only the elephants.

The remainder of the paper is organized as follows. In Section 2, we present the theory and online framework for classifying traffic flows into elephants and mice. This is followed by Section 3 which briefly presents the theory for estimating the distribution of the elephants. Evaluating the effectiveness of our approach is carried out in Section 4 with conclusions in Section 5.

2 An Online Framework for Identifying the Elephants

In this work, we define traffic *flows* to refer to packets with similar attributes. For example, a flow might be defined to consist of packets having identical values of five-tuple (source address, destination address, source port, destination port, protocol) or might be defined to comprise of packets matching specific payload information (e.g., group of all TCP packets with payload containing the string “crewman”). Thus, flows can be characterized by packet headers, payloads or a combination of both. The *size* of a flow is the number of packets (or bytes) belonging to the flow and the *duration* of a flow is its lifetime. Let $[\mathcal{F}] = \{F_1, F_2, \dots, F_i, \dots, F_N\}$ be a sequence of N FlowIDs $\{1, 2, \dots, i, \dots, N\}$, where each FlowID, F_i , is an index (i.e., a number between 1 and N) used to identify each flow in the underlying traffic. Denote $|F_i|$ to represent the number of packets belonging to the flow with FlowID F_i . It is important to note that the sequence $[\mathcal{F}]$ is sorted by increasing cardinality of the number of packets present in each FlowID. Under such circumstances, the flow classification problem is to identify and separate the F_i s that define the elephants and the mice. Now let us now consider an ergodic and discrete random process where each F_i is an independent variable drawn from the state space of $[\mathcal{F}]$. The state space of $[\mathcal{F}]$ consists of all possible FlowIDs. However, the random variables are not identically distributed. Denote $\{f_i\}$ to be the set of possible outcomes of F_i with $f \in [\mathcal{F}]$. Let us represent the probability mass function (pmf) of the sequence $\{F_i\}_{i=1}^N$ by: $P(F_1 = f_1, \dots, F_N = f_N) = p(f_1, \dots, f_N)$ Let $H(\mathcal{F}) = H(F_1, F_2, \dots, F_N)$ denote the *joint entropy* of the sequence $\{F_i\}_{i=1}^N$ and denote $\bar{H}_{\mathcal{F}}$ to be the *entropy rate* of $\{F_i\}_{i=1}^N$. Then, $H(\mathcal{F})$ and $\bar{H}_{\mathcal{F}}$ are defined as follows [18]:

$$H(\mathcal{F}) = H(F_1, F_2, \dots, F_N) = \sum_{i=1}^N H(F_i | F_{i-1}, \dots, F_1) \quad (1)$$

$$\bar{H}_{\mathcal{F}} = \frac{1}{N}H(\mathcal{F}) \quad (2)$$

Since according to our assumption, the F_i s are independent, Equation (1) reduces to: $H(\mathcal{F}) = \sum_{i=1}^N H(F_i)$ which is the summation of the individual entropies of the flow. At this point, it is worth mentioning that it is possible to *estimate* $H(\mathcal{F})$ without considering individual flow entropies [13]. However, this is not considered in this work.

Definition 1. *The set of elephants present in a sampled traffic is represented by the sequence, $\{F_1 F_2 \dots F_{N'}\}$, where $N' \ll N$ denotes the total number of elephants.*

Definition 1 provides us with the set of all packets which belong to the set of elephants. Since our aim is to identify the sequence $\{F_1, F_2, \dots, F_{N'}\}$, we need to isolate the sequence of FlowIDs that form the high probability set. If we visualize the set, $[\mathcal{F}]$, as an information source, then the existence of the above sequence of FlowIDs is governed by the probability of occurrence of a jointly typical sequence based on AEP. Note that the results based on AEP hold true only when the number of FlowIDs present in the sampled traffic volume is very large. Now considering the fact that there can be several sets of typical sequences, we have the following lemma for the set of elephants:

Lemma 1. *For traffic volumes with large number of FlowIDs (i.e., $N \rightarrow \infty$), the occurrence of the sequence $\{F_1, F_2, \dots, F_{N'}\}$, $N' \ll N$, is equiprobable and approximately equal to $2^{-N\bar{H}_{\mathcal{F}}}$.*

Lemma 1 follows directly from the property of AEP. In view of the above, we can say that out of all the possible FlowIDs, that sequence which belongs to the typical set has the maximum concentration of probability. The sequences outside the typical set are *atypical* and their probability of occurrence is extremely low. As evident from the above lemma, a typical sequence implies that FlowIDs in the typical set are associated with a large number of packets. If we consider the distribution of FlowIDs in the Internet traffic, we can easily correlate this property with the Zipf distribution of Internet flows. Hence, it is not surprising that most of the elephants belong to the typical set. However, what is the guarantee that such a sequence really exists?

Definition 2. *The joint entropy, $H(\mathcal{F})$ for a stationary, stochastic process is a decreasing sequence in N and has a limit equal to its entropy rate.*

Definition 2 implies that the probability of correctly identifying elephants *increases* with the corresponding increase in traffic volume. This observation is of fundamental nature since it enables us to *scalably* create an approximate list of LLFs (i.e., a typical sequence), while avoiding needless complex computation.

2.1 Algorithm for Flow Classification

Let $L \{\}$ be the list of empty LLFs and m the number of FlowIDs observed at the time instant the classification algorithm is being executed. Denote P_i^p and P_i^b to indicate the probability of occurrence of FlowID F_i in the sequence F_1, \dots, F_N , when considering the number of packets and payload bytes, respectively. Then,

$$P_i^p = |F_i| / \sum_{i=1}^n |F_i| \text{ and } P_i^b = \frac{\text{(cumulative payload carried by } F_i)}{\text{(total bytes observed at time instant t)}} \quad (3)$$

The pseudo-code of the classification algorithm is as follows:

1. Initialize list $L \{\} := \text{null}$
2. $m :=$ number of FlowIDs in current context;
- Loop: over all sampled $\{F_i\}$
 3. calculate probability P_i (Note: if the aim is to identify the set of elephants based on the number of packets, replace P_i with P_i^p . Similarly, for identifying the set of elephants based on payload size, replace P_i with P_i^b), for each $\{F_i\}$ using Equation 3
 4. calculate $H(\mathcal{F})$ and $\bar{H}_{\mathcal{F}}$
 5. if $p(F_i) \geq 2^{-n\bar{H}_{\mathcal{F}}}$
 6. add F_i to L
- Done
7. List $L \{\}$ contains the set of traffic flows which are elephants.

3 Estimating the Density of Elephants Flows

We employ the Parzen window [19] technique (explained below) on the set \mathcal{F} for determining the density of the identified elephants. Note that the likelihood estimator from the coupon collector problem [23] can be employed on the sampled set of all elephants in order to identify the set of all elephants present in the underlying traffic. However, that aspect is not presented in this work. The standard method is to choose a well-defined kernel function (e.g., Gaussian) of definite width and convolve it with the known data points. Let $\hat{f}_h(x)$ be the pdf of the random variable \mathcal{X} we are trying to estimate for the set \mathcal{F} and be defined as [19]:

$$\hat{f}_h(x) \approx \frac{1}{Nh} \sum_{i=1}^N \psi\left(\frac{x - x_i}{h}\right) \quad (4)$$

where $\{x_i\}_{i=1}^N$ are the data points of \mathcal{X} and $\psi(\cdot)$ is a suitable kernel smoothing function of width h , also referred to as the bandwidth of $\psi(\cdot)$. In this approach, the estimated pdf is a linear combination of kernel functions centered on individual x_i . In Equation (4), the bandwidth factor h is the most important term in the estimation process [20]. The optimal value of the kernel window h can

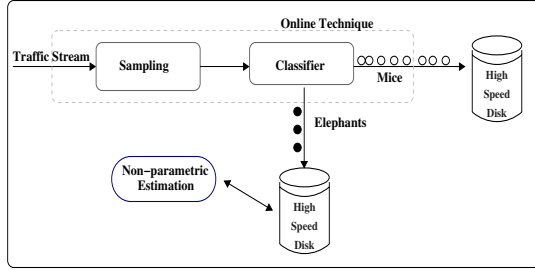


Fig. 1. On-line classification of traffic streams. The pdf of elephant flows in sampled traffic stream is estimated using the non-parametric Parzen window technique.

be calculated by minimizing the integrated mean square error (IMSE) between $f(x)$ (original pdf) and $\hat{f}_h(x)$; i.e.,

$$\text{minimize } \left\{ \int \left\{ \hat{f}_h(x) - f(x) \right\}^2 dx \right\}.$$

In general, the process of finding the optimal window size is cumbersome as we do not know beforehand the nature of the density function that we are trying to estimate. Since the shape (degree of smoothness) of $\hat{f}_h(x)$ is closely related to the kernel function used, we use the Gaussian kernel function to eliminate “noises” in the pdf estimation. Thus:

$$\psi(u) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right) \quad (5)$$

Corresponding to the Gaussian kernel, the bandwidth h can be approximated using *Silverman’s rule of thumb* [21] that satisfies the IMSE criteria. Consequently, h is defined as: $h = 1.06 \hat{\sigma} N^{-1/5}$ where $\hat{\sigma} = \sqrt{\frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N}}$ denotes the standard deviation of the sample.

4 Performance Evaluation

In this section, we evaluate the performance of our algorithm using packet traces obtained from NLNR [15]. We compare our approach with the results of Mori et al. [8] in the figures) for comparing the number of elephants detected in the traffic stream. Specifically, we use three traces: (i) 20040130-133500-0.gz, (ii) 20040130-13400-0.gz, and (iii) 20040130-134500-0.gz. The cumulative duration of the three files is 900 seconds and contains 23.2 million packets. They subsequently map to 618, 225 FlowIDs, where each FlowID is defined using the number of packets.

4.1 Identifying the Elephants

In Figure 2, we plot the number of elephants predicted using our classification algorithm and compare it with the approach of [8]. We have used the frequency

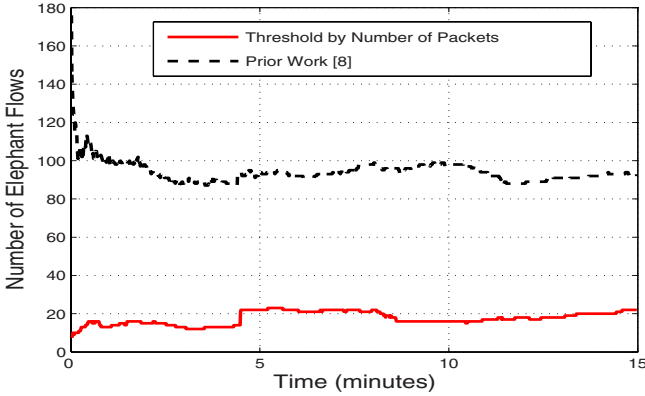


Fig. 2. Time series of the occurrence of elephants as estimated using our algorithm versus as predicted using the approach of [8]

of occurrence of the packets as the basis for calculating the flow probabilities. Apart from the already known facts that the proportion of elephants are small in number (0.0035% in our case), two important conclusions can be immediately drawn from this figure:

- the set of elephants detected during the *initial phase* (first 5 minutes for the traffic traces under consideration) of our algorithm identifies FlowIDs that exhibit *bursty behavior*. On close analysis of the traffic traces, we found that this is indeed the case and is due to the fact that such FlowIDs cause immediate concentration of the probability mass function of the entire traffic sample.
- the *proportion of elephants* classified using the frequency of occurrence of packets (i.e. probability P_i^p) is almost equal in extent to those detected by considering the volume (bytes) of traffic (i.e. probability P_i^b). Notice that, using the approach of [8], the number of elephants are estimated at around 85 – 90. If the traffic traces beyond 5 minutes are considered (not shown in this study), the approach of [8] exhibits a *decreasing trend*.

The occurrence of mice, however, shows well-established behavior. They are large in number but grow with the continuation of traffic stream.

4.2 Traffic Distribution: Elephants and Mice

In Figures 3 and 4, we analyze the traffic carried by elephants and mice when considering the frequency of occurrence of packets and the the volume of traffic carried by each FlowID as the basis of flow probability calculation. While 99% FlowIDs carry 70% of network traffic, elephant flows (less than 1%) carry 30% of the traffic. Such dynamics is unaffected if we choose the frequency or the volume of traffic as the basic for probability calculation.

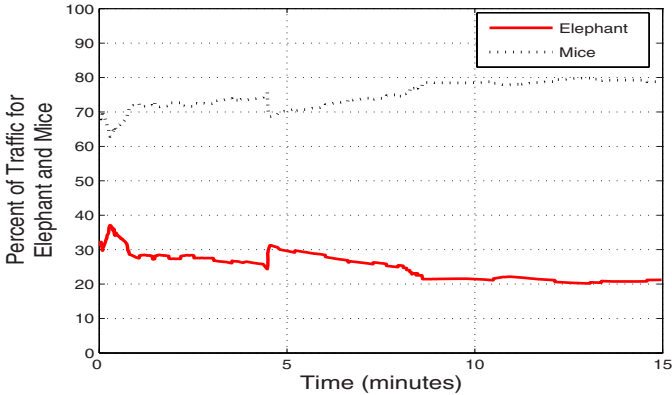


Fig. 3. Traffic Distribution of Elephant Flows considering the frequency of occurrence of packets in each flow (i.e., P_i^p)

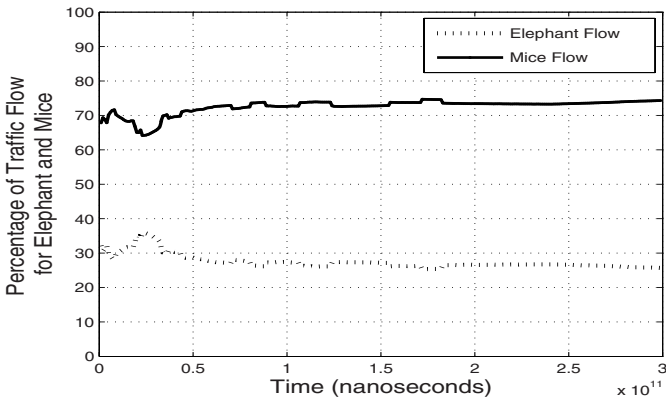


Fig. 4. Traffic Distribution of Elephant Flows considering the volume of traffic in each flow (i.e., P_i^b)

4.3 Entropy Distribution: Elephants and Mice

In Figures 5 and 6, we show the *temporal variation* of the entropy of the ratio of entropy between the elephants and mice, where only of the FlowIDs considered to be mice. During the first 500msecs of input traffic, we observe a *dip* in the entropy of the elephant flows. This is due to the presence of bursty elephant flows which causes a temporary increase in the probability of the set of elephants. However, as the experiment continues, the entropy of the mice increases while the entropy of the elephant flows decreases. Since the entropy of the typical set is a *decreasing sequence* with respect to the number of FlowIDs, the probability and proportion of FlowIDs classified as elephants increases.

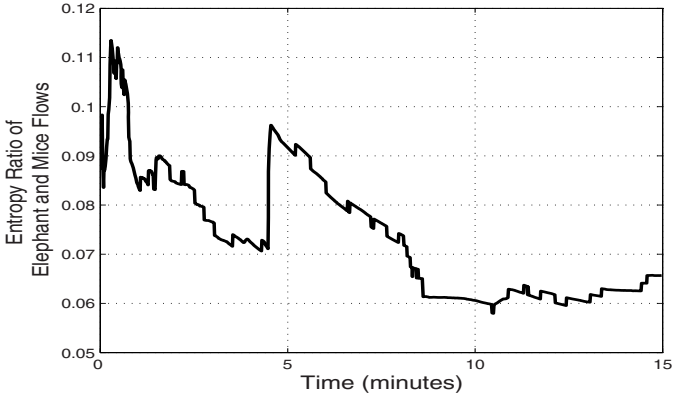


Fig. 5. Ratio of entropy between the elephants and mice

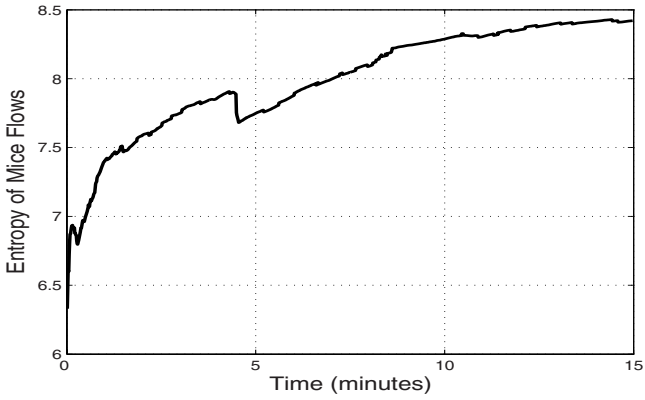


Fig. 6. Temporal distribution of Entropy of mice

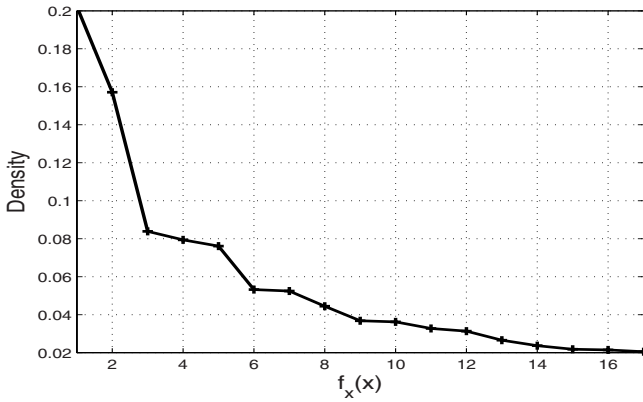


Fig. 7. Original pdf of the elephants present in the traffic stream

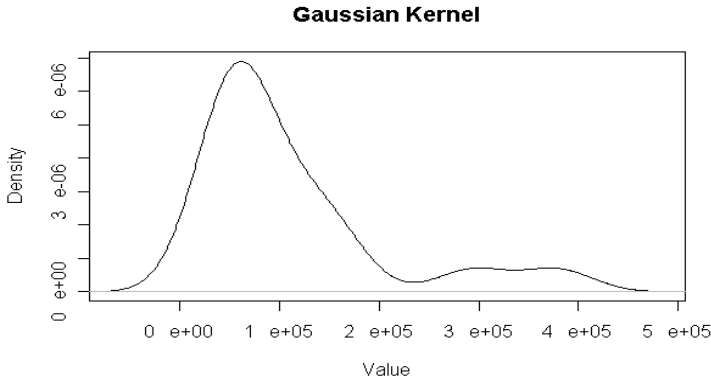


Fig. 8. Pdf estimated using parzen window technique

This unique trend of entropy variation guarantees conservative, yet accurate flow classification of high traffic volumes.

Estimating the density function of the distribution. In Figures 7 and 8, we plot the pdf of the elephants present in the original traffic stream versus the pdf of elephants identified using our approach based on the Parzen window technique. Observe that the trend observed is similar in both the cases.

5 Conclusions

In this paper we have focused on classifying and estimating the properties of elephants and mice based on AEP from Information Theory. Although considerable attention has been directed in identifying high and low traffic volumes, we feel that our approach using typical sequences simplifies the problem to a large extent and provides a standard yardstick for defining such long-lived-flows. We have evaluated our algorithm with the approach of [8], and have observed that our approach is able to identify bursty elephant flows and at the same time does not overestimate the number of occurrence of the elephants. As part of future work, we would like to carry out these observations on NLANR trace of more than one hour duration and see how the classification and estimation algorithms perform if the input traffic becomes smooth with non-negligent coefficient of variation [14].

References

1. J. S. Marron, F. Hernandez-Campos and F. D. Smith, “*Mice and Elephants Visualization of Internet Traffic*”, available online at citeseer.ist.psu.edu/531734.html.
2. A. Kuzmanovic and E. Knightly, “*Low-Rate TCP-Targeted Denial of Service Attacks (The Shrew vs. the Mice and Elephants)*”, in *Proceedings of ACM SIGCOMM*, August 2003.

3. Y. Joo, V. Riberio, A. Feldmann, A. C. Gilbert and W. Willinger, "On the impact of variability on the buffer dynamics in IP networks", *Proc. 37th Annual Allerton Conference on Communication, Control and Computing*, September 1999.
4. K. Papagiannaki, N. Taft, S. Bhattacharyya, P. Thiran, K. Salamatian and C. Diot "On the Feasibility of Identifying Elephants in Internet Backbone Traffic", *Sprint ATL Technical Report TR01-ATL-110918*, November 2001.
5. N. Duffield, C. Lund and M. Thorup, "Properties and Prediction of Flow Statistics from Sampled Packet Streams", in *Proc. of ACM SIGCOMM Internet Measurement Workshop*, Nov. 2002.
6. N. Duffield, C. Lund and M. Thorup, "Estimating Flow Distributions from Sampled Flow Statistics", in *Proc. of ACM SIGMETRICS*, August 2003.
7. A. Kumar, J. Xu, O. Spatschek and L. Li, "Space-Code Bloom Filter for Efficient Per-Flow Traffic Measurement", *IEEE INFOCOM*, August 25-29, 2004.
8. T. Mori, M. Uchida, R. Kawahara, J. Pan and S. Goto, "Identifying elephant flows through periodically sampled packets", in *Proceedings of the ACM SIGCOMM Workshop on Internet Measurement Workshop (IMW)*, 2004.
9. J. Sommers and P. Barford, "Self-Configuring Network Traffic Generation", in *Proc. of ACM SIGCOMM Internet Measurement Workshop*, October 25-27, 2004.
10. C. Estan and G. Varghese, "New directions in traffic measurement and accounting: Focusing on the elephants, ignoring the mice", *ACM Trans. Comput. Syst.*, vol. 21, no. 3, 2003.
11. K. Papagiannaki, N. Taft, S. Bhattacharyya, P. Thiran, K. Salamatian and C. Diot, "A pragmatic definition of elephants in internet backbone traffic", in *Proceedings of the ACM SIGCOMM Workshop on Internet Measurement Workshop (IMW)*, 2002.
12. J. Wallerich, H. Dreger, A. Feldman, B. Krishnamurthy and W. Willinger, "A Methodology for Studying Persistency Aspects of Internet Flows", in *ACM SIGCOMM Computer Communication Review*, vol. 35, Issue 2, pp. 23 - 36, 2004.
13. A. Lall, V. Sekhar, M. Ogihara, J. Xu and H. Zhang, "Data Streaming Algorithms for Estimating Entropy of Network Traffic", in *Proc. of ACM SIGMETRICS*, June 2006.
14. J. Cao, W. S. Cleveland, D. Lin and D. X. Sun, "Internet Traffic Tends Towards Poisson and Independent as Load Increases", in *Nonlinear Estimation and Classification*, New York, Springer-Verlag, 2002.
15. NLANR AMP Website: <http://pma.nlanr.net/Special/>
16. N. Brownlee, "Understanding Internet Traffic Streams: Dragonflies and Tortoises", *IEEE Communications Magazine*, Cot. 2002.
17. A. Kumar, M. Sung, J. Xu and L. Wang, "Data Streaming Algorithms for Efficient and Accurate Estimation of Flow Size Distribution", *ACM SIGMETRICS*, August 25-29, 2003.
18. T. M. Cover and J. A. Thomas "Elements of Information Theory", John Wiley, 1991.
19. E. Parzen, "On estimation of a probability density function and mode," *Time Series Analysis Papers*. San Francisco, CA: Holden-Day, 1967.
20. S. Raudys, "On the effectiveness of Parzen window classifier," *Informatika*, vol. 2, no. 3, pp 435-454, 1991.
21. B. W. Silvean, *Density Estimation for Statistics and Data Analysis*, Chapman and Hall, 1986.
22. S. R. Kundu, B. Chakravarty, K. Basu, and S. K. Das, "FastFlow: A Framework for Accurate Characterization of Network Traffic", *IEEE ICDCS*, 2006.
23. M. Finkelstein, H. G. Tucker and J. A. Veeh, *Confidence Intervals for the number of Unseen Types*, *Statistics and Probability Letters*, vol. 37, pp. 423-430, 1998.