

Estimating Original Flow Length from Sampled Flow Statistics

LIU Weijiang^{1,2,3} GONG Jian^{2,3} DING Wei^{2,3} PENG Yanbing^{2,3}

¹ Post Doctoral Station for Computer Science and Technology, Southeast University,
NanJing, Jiangsu, China 210096

² Department of Computer Science and engineering, Southeast University, NanJing, Jiangsu,
China 210096

³ Jiangsu Province Key Laboratory of Computer Networking Technology, NanJing, Jiangsu,
China 210096

wjliu@njnet.edu.cn

Abstract. Packet sampling has become an attractive and scalable means to measure flow data on high-speed links. Passive traffic measurement increasingly employs sampling at the packet level and makes inferences from sampled network traffic. This paper proposes a maximum probability method that estimates the length of the corresponding original flow from the length of a sampled flow. We construct the probability models of the original flow length distributions of a sampled flow under the assumptions of various flow length distributions, respectively. Through recovery analyzing with different parameters, we obtain a consistent linear expression that reflects the relationship between the length of sampled flow and that of the corresponding original flow. Furthermore, after using publicly available traces and traces collected from CERNET to do recovery experiments and comparing the experiment outcomes and theoretic values calculated with Pareto distributions, we may conclude that the maximum probability method calculated by using the Pareto distribution with 1.0 can be used to estimate original flow length in the concerned network.

Key words: Measurement, Packet sampling, Estimation, IP Flows, Conditional Probability

1 Introduction

Internet is the largest global internetwork, connecting hundreds of million of computers worldwide. With the rapid increase of network service types and user number, **measuring** the volume of network traffic and network performance is becoming more **difficult**. MCI **measuring traces** showed that there were more than 250,000 flows in 1997[1]. Now a minute of traffic in high-speed links (40Gbps) can contain millions of flows. The increasing speed of network links makes it infeasible to collect complete data on all packets or network flows. This is due to the costs and scale of the resources required to accommodate the data in the measurement infrastructure. These constraints motivate reduction of the data. Packet sampling

allows the retention of arbitrary detail while at the time reducing data volumes, so that it has become an attractive and scalable means to measure flow data. In 1993, K. Claffy et al [2] studied systematic, stratified random and simple random sampling method by packet or time. The method of sampling by packet's content was discussed in [3]. The Packet Sampling working group (PSAMP), which was founded by IETF in 2003, is chartered to define a standard set of capabilities for network elements to sample subsets of packets by statistical and other methods [4]. The method and application of Trajectory Sampling was considered in [5][6]. Furthermore, sampling techniques have been employed in network products such as Cisco's Netflow [7] and NetranMet [8].

Sampling entails an inherent loss of information. We wish to use statistic inference to recover information as much as possible. The original flow length is very important for many applications. This paper proposes a Maximum Probability (MP) method that estimates the length of the corresponding original flow from the length of a sampled flow. This method is simple to calculate and easy to implement.

The main contributions of this paper include:

- 1) Maximum Probability (MP) method that estimates the length of the corresponding original flow from the length of a sampled flow is proposed.
- 2) A linear expression that reflects the relationship between the length of sampled flow and that of the corresponding original flow is obtained. Although it is different for different Pareto parameters, the difference is very small.
- 3) **we conclude** that the value calculated by using the Pareto distribution with 1.0 can be used to estimate original flow length In the concerned network.

The rest of this paper is organized as follows. In the next section, we review some elementary concepts on flow and sampling. In section 3 we construct the probability models of the original flow length distributions of a sampled flow under the assumptions of various flow length distributions, respectively. In section 4, we propose the Maximum Probability (MP) method and compare the estimating results of uniform distribution, Pareto distributions and empirical distributions. In section 5, we compare some related work. We conclude with some proposals for future work in section 6.

2 Some Elementary Concepts

This paper considers sampling some target proportion $p=1/N$ of the packet stream. There are a number of different ways to implement this, e.g. independent sampling of packets with probability $p=1/N$, and periodic selection of every N^{th} packet from the full packet stream. In both cases we will call N the sampling period, i.e., the reciprocal of the average sampling rate. Although the length distributions by random and periodic sampling can be distinguished, the differences are, in fact, sufficiently small [9]. There are at least a few definitions for the term flow depending on the context of research. In this study, we employ the one adopted in [10] which stems from the packet train model by Jain and Routhier [11].

Definition 1. A flow is defined as a stream of packets subject to flow

specification and timeout.

When a packet arrives, the specific rules of flow specification determine which active flow this packet belongs to, or if no active flow is found that matches the description of this packet, a new flow is created. On the other hand, flow timeout determines when to terminate a flow. When no new packet arrives within a given flow timeout period since the last packet arrived, this flow is terminated and the arrival time of the last packet becomes the end point of this flow. In this paper, the flow interpacket timeout is 64 seconds. A TCP flow is a stream of TCP packets subject to timeout and having the same source and destination IP addresses, same source and destination port numbers. Similarly, a UDP flow is a stream of UDP packets associated with above specification. A general flow is a stream of packets subject to timeout and having the same source and destination IP addresses, same source and destination port numbers (not considering protocol). In this paper, we will use the term original flow to describe the above flow. A flow length is the number of packets in the flow.

Definition 2. A sampled flow is defined as a stream of packets that are sampled at probability $p=1/N$ from an original flow.

3 Probability Distribution of Original flow length

In this paper, sampling probability is $p=1/N$. For a **specific** original flow F , let X_F denote the number of packets in F , Y_F denote the number of packets in the sampled flow from F . The conditional distribution of Y_F , given that $X_F=l$, follows a binomial

distribution $\Pr[Y_F=k | X_F=l] = B_p(l, k) = \binom{l}{k} p^k (1-p)^{l-k}$. For an original flow F ,

let $\Pr[Y_F=y, X_F=x]$ denote the probability of $X_F=x$ and $Y_F=y$, by the conditional probability formula,

$$\Pr[X_F = x | Y_F = y] = \frac{\Pr[Y_F = y, X_F = x]}{\Pr[Y_F = y]} = \frac{\Pr[Y_F = y | X_F = x] \Pr[X_F = x]}{\Pr[Y_F = y]} \quad (1)$$

and by the complete probability formula,

$$\begin{aligned} \Pr[Y_F = y] &= \sum_{i=y}^{\infty} \Pr[Y_F = y | X_F = i] \Pr[X_F = i] \\ &= \sum_{i=y}^{\infty} B_p(i, y) \Pr[X_F = i], \quad y=0,1,\dots \end{aligned} \quad (2)$$

3.1 Uniform Distribution

Suppose original flows lengths satisfy uniform distribution, that is , $\Pr[X_F=k]=\Pr[X_F=k+1]$, for all $k=1,2, \dots$. Moreover,

$$\begin{aligned} \sum_{l=k}^{\infty} B_p(l,k) &= \sum_{l=k}^{\infty} \binom{l}{k} p^k (1-p)^{l-k} = p^k \sum_{l=0}^{\infty} \binom{l+k}{k} (1-p)^l \\ &= p^k \sum_{l=0}^{\infty} \binom{l+k}{k} q^l = p^k (1-q)^{-k-1} = 1/p = N \end{aligned}$$

hence $\Pr[Y_F = y] = \sum_{i=y}^{\infty} B_p(i,y) \Pr[X_F = i] = \Pr[X_F = y] \sum_{i=y}^{\infty} B_p(i,y) = \Pr[X_F = y]/p$, so

$$\Pr[X_F = x | Y_F = y] = \frac{\Pr[Y_F = y | X_F = x] \Pr[X_F = x]}{\Pr[Y_F = y]} = \frac{\Pr[Y_F = y | X_F = x] \Pr[X_F = x]}{\Pr[X_F = y]/p} = p B_p(x, y)$$

and we obtain the following results:

Lemma 1 The probability that a sampled flow of length k is sampled from an original flow of length l is $\Pr[X_F = l | Y_F = k] = \binom{l}{k} p^{k+1} (1-p)^{l-k}$, $l=k, k+1, \dots$.

Lemma 2 The mean and variance of the above probability distribution are $EX = N(k+1) - 1$ and $DX = (N+1)N(k+1)$, respectively.

Let $a_1 = \frac{B_p(l,k)}{B_p(l-1,k)} = \frac{(l-1)(1-p)}{l-k-1} = 1 + \frac{k-(l-1)p}{l-k-1}$. For $l < kN+1$, since $a_1 > 1$, hence $B_p(l,k)$ is increasing as l increases. For $l > kN+1$, since $a_1 < 1$, hence $B_p(l,k)$ is increasing as l decreases. At $l=kN+1$, $a_1 = 1$, $B_p(l,k) = B_p(l-1,k)$ is maximized.

Lemma 3. The probability $\Pr[X_F = l | Y_F = k]$ is maximized at $l=kN, kN+1$. It is increasing as l increases for $l < kN+1$ and decreasing as l increases for $l > kN+1$.

3.2 Pareto Distribution

Assume original flow lengths satisfy Pareto distribution. The probability density function is defined as

$$\Pr[X_F=x] = ba^b / x^{b+1}, \quad x=1,2,\dots \quad (3)$$

Where β is called Pareto parameter. Hence formula (2) can be written as:

$$\Pr[Y_F = y] = \sum_{i=y}^{\infty} B_p(i, y) \Pr[X_F = i] = \sum_{i=y}^{\infty} B_p(i, y) b a^b / i^{b+1}, y=0,1,\dots$$

Lemma 4. Under the assumption that original flow lengths satisfy Pareto distribution, the probability that a sampled flow of length y is sampled from an original flow of length x is

$$\Pr[X_F = x | Y_F = y] = \frac{B_p(x, y) / x^{b+1}}{\sum_{i=y}^{\infty} B_p(i, y) / i^{b+1}}.$$

Lemma 5. The mean of the above probability distribution, given $Y_F=y$, ($y \neq 0$), is

$$E[X_F | Y_F = y] = \sum_{x=y}^{\infty} x \Pr[X_F = x | Y_F = y] = \frac{\sum_{i=y}^{\infty} B_p(i, y) / i^b}{\sum_{i=y}^{\infty} B_p(i, y) / i^{b+1}}.$$

4 Maximum Probability Method (MP)

The purpose of the Maximum Probability method (MP) is to estimate the length of the corresponding original flow from the length of a sampled flow. The point we are trying to make is that MP estimates the length of the original flow according to maximum probability. MP contains three steps.

i) Computing probability. Given a sampled flow with fixed length k , compute $\Pr[X_F=l | Y_F=k]$ for $l=k, k+1, \dots$.

ii) Finding maximum probability. Define $mp = \max_{l \geq k} \{\Pr[X_F = l | Y_F = k]\}$.

iii) Estimating length. Let $\hat{l} = \min_l \{l | mp = \Pr[X_F = l | Y_F = k]\}$, we write

our estimate of the length of the original flow as \hat{l} .

Below we apply the MP method to different flow length distributions.

4.1 Uniform Distribution

Let the length of original flows be uniform distribution. By lemma 3, we can use MP to obtain a linear expression as

$$\overline{X}_F = \frac{Y_F}{p} \quad (4)$$

where \overline{X}_F is the estimate of the length of the original flow, Y_F is the length of the sampled flow, p is sampling probability.

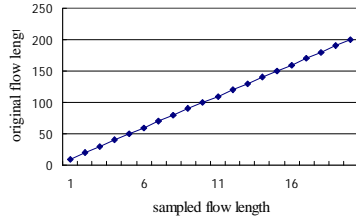


Figure 1: MP estimate of original flow length under uniform distribution

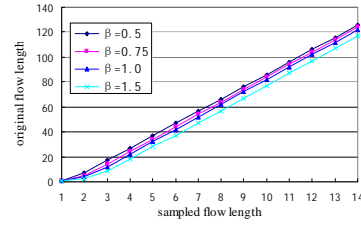


Figure 2: MP estimate of original flow length under Pareto distribution

The estimates of the original flow lengths for $p=0.1$ are shown in figure 1. From this figure we can observe the linear relationship between the length of original flow and that of sampled flow under uniform distribution. 其它概率情况呢? 为什么只选 $p=0.1$?

4.2 Pareto Distribution

Let the length of original flows be Pareto distribution. By lemma 4, we can compute the results with $\beta = 0.5, 0.75, 1.0, 1.5$, respectively, and obtain an expression that reflect the relationship between the estimated length of the original flow and the sampled flow length as

$$\overline{X}_F = \frac{Y_F}{p} - n(p, \beta) \quad \text{for } Y_F \geq 1/p \quad (5)$$

where \overline{X}_F is the estimate of the length of the original flow, Y_F is the length of the sampled flow, p is sampling probability, $n(p, \beta)$ is a positive integer involving p and β as defined in following section and 下节说的是同一个东西吗?. In figure 2, with sampling probability $p=0.1$, and Pareto parameter $\beta = 0.5, 0.75, 1.0, 1.5$, respectively, we can observe that, for a sampled flow with fixed length Y_F , \overline{X}_F is increasing as β decreases. It is minimized at $\beta = 1.5$, maximized at $\beta = 0.5$, medial at $\beta = 1.0$. Though there are differences for different parameter, the estimates show very similar tendencies for all parameters. Therefore, we may conclude that \overline{X}_F calculated by using the parameter $\beta = 1.0$ can be used as approximations under unknown parameter values in the concerned network. 同

样，为什么只选 $p=0.1$?

4.3 Finite Pareto Distribution

In the concerned network, due to constraint of measurement time, the lengths and numbers of flows all are finite. Let M denote the maximum original flow length in an original flow distribution, the probability density function in formula (3) is

$$\Pr[X_F = x] = g / x^{b+1}, \quad x=1, \dots, M \quad (6)$$

where $g = \frac{ba^b}{\sum_{x=1}^M ba^b / x^{b+1}}$, $0 < g < 1$. Formula (6) can be extended so that β

can be extended to real field, that is, β can be zero and negative. We call formula (6) extended Pareto distribution.

For $\beta = 0$, formula (6) is

$$\Pr[X_F = x] = g / x, \quad x=1, \dots, M \quad (7)$$

For $\beta = -1$, formula (6) is

$$\Pr[X_F = x] = g, \quad x=1, \dots, M \quad (8)$$

Formula (8) is uniform distribution, therefore we call uniform distribution as degenerate Pareto distribution.

We now consider how the finite constraint impacts the conditional probability. Due to the constraint of finite number of flows, the conditional probability in lemma 4 is written by

$$\Pr_M[X_F = x | Y_F = y] = \frac{B_p(x, y) / x^{b+1}}{\sum_{i=y}^M B_p(i, y) / i^{b+1}} \quad (9)$$

Obviously $\Pr_M[X_F = x | Y_F = y] = r(y) \Pr[X_F = x | Y_F = y]$, where

$$r(y) = \frac{\sum_{i=y}^M B_p(i, y) / i^{b+1}}{\sum_{i=y}^{\infty} B_p(i, y) / i^{b+1}} < 1 \text{ is a function with variable } y. \text{ Hence, for a fixed } y,$$

the above two probabilities are maximized at the same x . Therefore formulae (4) and (5) are still valid, we rewrite them as consistent form

$$\overline{X}_F = \frac{Y_F}{p} - n(p, \beta) \quad \text{for } Y_F \geq 1/p \quad (10)$$

where \overline{X}_F is the estimate of the length of the original flow, Y_F is the length of the sampled flow, p is sampling probability, $n(p, \beta)$ is a binary function with variables p and β whose value domain is integer set. Function $n(p, \beta)$ has the following properties:

- i) It is a monotone decreasing function on variable p , that is, for fixed β , is decreasing as p increases.
- ii) It is a monotone increasing on variable β , that is, for fixed p , is increasing as p increases.
- iii) $n(p, -1) = 0$, for any p ($0 < p < 1$).

For example, $n(0.1, -0.5) = 5$, $n(0.1, 0) = 10$, $n(0.1, 0.5) = 14$, $n(0.1, 0.75) = 16$, $n(0.1, 1.0) = 18$, $n(0.1, 1.25) = 21$, $n(0.1, 1.5) = 23$, $n(0.1, 2.0) = 27$, $n(0.1, 3.0) = 36$.

4.4 Comparing for Different Distributions

We use six traces to **verify the MP method**. The first three traces [10], **all containing** packets during a 10 minute period, were collected with a [Dag3.2E](#) 10/100 MBit/sec Ethernet card at the outside of the firewall servicing researchers at Bell Labs via a 9 MBits/sec link to the Internet in May 2002. Gateway link to ISP is configured at 9 MBps. Local network hosts serve 400 people, mostly technical, and about 50 administrative staff. All flows are bidirectional, routing is fully symmetric. The other three Traces, either of which contains packets during a 10-minute period too, were collected at Jiangsu **provincial** network border of **China Education and Research Network (CERNET)** in disjoint time interval on April 17, 2004. The backbone **capacity** is 1000Mbps; mean traffic per day is 587 Mbps. For each trace, we sample at $p=0.1, 0.05, 0.01$ respectively. Then we use MP to estimate the length of original flow. We find the estimated lengths are very close at same sampling probability for the six traces. For clear display, we only show the estimates in three experiments at sampling probability $p=0.1$ in figure 3. As shown in figure 3, the estimates are very close.

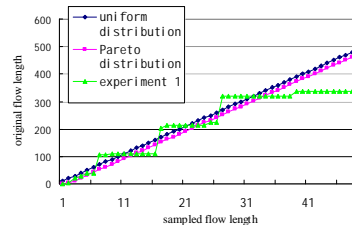
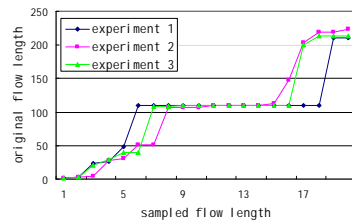


Figure 3: MP estimate of original flow length Figure 4: Comparing MP estimates of original

under empirical distribution

flow length under different distributions

Figure 4 illustrates the estimates for uniform distribution, Pareto distribution ($\beta = 1.0$) and experiment 1 实验的定义是什么? . We can observe that the estimates for uniform are slightly big, but the estimates for Pareto distribution with $\beta = 1.0$ fit those for experiment 1 very well. From figure 5, we can see that the estimate results of several experiments all move up and down at the point of estimate for Pareto distribution ($\beta = 1.0$) . Therefore we can use the Pareto distribution with $\beta = 1.0$ as theoretical distribution of the three empirical distributions.

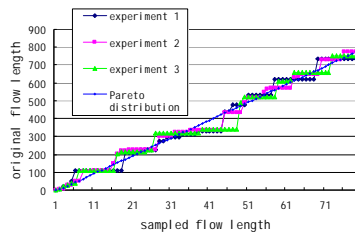


Figure 5: Comparing MP estimates of original flow length under empirical distributions and Pareto distribution

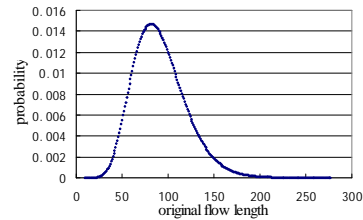


Figure 6: Probability distribution of original flow length, given length 10 of sampled flow at sampling probability $p=0.1$

4.5 Difficulties and Applications

When we use MP to estimate the original flow length of a sampled flow, we must find the **one** that makes the probability maximized. **However**, the maximized probability sometimes is **still** very small. For example, consider a sampled flow with length 10 at sampling probability 0.1. We calculate the probability by using Pareto distribution with 1.0 to estimate its original flow length. Figure 6 displays the probability distribution of the original flow length of the sampled flow. **Within it** we can observe that the probability at length 82 is maximized with value 0.0147, that is, the probability that the sampled flow is sampled from the original flow of length 82 is 1.47%. Although 0.0147 is maximum value, it is too small. It reflects the difficulty and uncertainty associated with an estimate. To improve certainty, we can estimate the confidence interval for original flow length subject to interval width (as small as possible). Suppose that the estimated length may fall into an interval with width 70, we can sum for some values and obtain a maximum probability $\Pr[51 \leq x \leq 120 \mid y=10]=0.800017$, that means that the estimated length will fall into the interval $[51,120]$ more than 80% of the time. Equivalently, $\Pr[61 \leq x \mid y=10]=0.88$ means that the estimated length will fall into the interval $[61, \infty]$ 88% of the time. **Therefore, we can at least use MP to estimate the boundary of length for a specific flow. 是吗?**

5 Related Work

Duffield et al, who started with study of obtaining the original flow information from the sampled flow statistics in [13], presented the idea of inferring the original flow statistics from the sampled flow statistics and showed how to infer the mean length of flows. However, since they fail to consider original flow distributions, the results are rough and hard to use. This is followed by [9] in which the flow distribution is inferred from the sampled statistics. After showing that the naive scaling of the flow distribution estimated from the sampled traffic is in general not accurate, the authors propose an EM algorithm to iteratively compute a more accurate estimation. Scaling method is simple, but it exploits the sampling properties of SYN flows to estimate TCP flow frequencies; EM algorithm does not rely on the properties of SYN flows and hence is not restricted to TCP traffic, but its versatility comes at the cost of computational complexity and its termination criterion is absent.

In 2004 Tatsuya Mori et al [12] developed techniques and schemes to identify elephant flows in periodically sampled packets, that is, to determine whether the original flow is an elephant flow by the sampled flow length. The key is to find the threshold of per-flow packets in sampled packets, which can reliably indicate whether or not a flow is actually an elephant flow in unsampled packets. The advantage of this approach is due to its simplicity. Periodic sampling without per-packet processing can be easily implemented. Its disadvantage is very hard to implement a proper trade-off of false positives and false negatives.

In 2005 Noriaki Kamiyama [13] proposed a short timeout method, which identifies high-rate flows with high accuracy from sampled packet information. This method is easy to implement and requires only a small amount of memory. In fact, high-rate flows are elephant flows in fixed time interval; so that identifying high-rate flow can be regard as identify some class of elephant flows. Our method is useful to not only identify elephant flows, but also estimate flow length distributions. *要进一步说明别人方法在流长度估计方面与你的差距/差别?*

6 ~~Conclusions and Future Work~~

This paper proposes a native method (MP) to estimates original flow length from the sampled flow. For different Pareto parameters, we obtain a consistent linear expression that reflects the relationship between the length of sampled flow and that of the corresponding original flow. The existence of the expression makes it very simple, easy and convenient to estimate original flow length. This expression also shows that the estimated results are very similar for various distributions. In the concerned network, the length distributions of flows collected in any time interval do not satisfy Pareto distributions with fixed parameter strictly, but they can follow a Pareto distribution with parameter in interval [0.5, 1.5] approximately. The Value 1.0 is the middle value of interval [0.5,1.5] exactly. Theory analysis and experiment

results show that it is a reasonable choice using parameter 1.0 to calculate at the condition of unknown parameter value.

7 Acknowledgements

This work is supported in part by the National Grand Fundamental Research 973 Program of China under Grant No.2003CB314804, the National High Technology Research and Development Program of China (2005AA103001), the Key Project of Chinese Ministry of Education under Grant No.105084, the Jiangsu Provincial Key Laboratory of Computer Network Technology No. BM2003201 and Jiangsu Planned Projects for Postdoctoral Research Funds.

References

1. Thomsom, k. Miller ,G.J., and Wilder,R.. Wide-area traffic patterns and characteristics[J]. IEEE Network Magazine, 11(6):10--23, November 1997 .
2. Claffy, K., Polyzos, G., Braun, H.: Application of Sampling Methodologies to Network Traffic Characterization[C]. May 1993, Proceedings of ACM SIGCOMM '93.
3. Cozzani, I., Giordano, S. A passive test and measurement system: traffic sampling for QoS evaluation[C]. Global Telecommunications Conference, 1998. GLOBECOM 1998. The Bridge to Global Integration. IEEE , Volume: 2 , 1998, Page(s): 1236 –1241
4. Packet Sampling (psamp), <http://www.ietf.org/html.charters/psamp-charter.html>, 2005-02-02.
5. Duffield, N. G., and Grossglauser, M.. Trajectory Sampling for Direct Traffic Observation[J]. IEEE/ACM Trans. on Networking, 9(3), 280-292, June 2001.
6. Duffield, N. G., and Grossglauser, M.. Trajectory Sampling with Unreliable Reporting[C]. IEEE Infocom 2004, March 7-11 2004 in Hongkong.
7. Sampled Cisco, http://www.cisco.com/en/US/products/sw/iosswrel/ps1829/products_feature_guide09186a_0080081201.html, 2002.12
8. NeTraMet Version 4.4 Now Available, <http://www2.auckland.ac.nz/net/Accounting/ntm.Release.note.html>, 2002.12.
9. Duffield, N.G, Lund, C. , Thorup, M.. Estimating Flow Distributions from Sampled Flow Statistics[C]. ACM SIGCOMM . 2003, Karlsruhe,Germany. August 25-29. 325-336.
- 10.Claffy, K.C., Braun,H.W., and Polyzos, G.C.: A parameterizable methodology for Internet traffic flow profiling. IEEE JSAC,13:1481-1494,1995.
- 11.Jain,R. and Routhier, S.A.: Packet trains-measurements and a new model for computer network traffic. IEEE JSAC, 4:986-995,1986.
- 12.NLANR: Abilene-I data set, <http://pma.nlanr.net/Traces/long/bell1.html>.

13. Duffield, N.G., Lund, C. , Thorup, M.. Properties and Prediction of Flow Statistics from Sampled Packet Streams[C]. ACM SIGCOMM Internet Measurement Workshop 2002, Marseille, France, November 6-8, 2002.
14. Tatsuya Mori, Masato Uchida, Ryoichi Kawahara. Identifying Elephant Flows Through Periodically Sampled Packets[C]. ACM SIGCOMM Internet Measurement Conference, 2004, Taormina, Sicily, Italy.
15. Noriaki Kamiyama. Identifying High-Rate Flows With Less Memory[C]. IEEE Infocom 2005, March 13-15, 2005 in Miami.