

## 高速ネットワークのためのインライン計測手法

Cao Le Thanh Man<sup>†</sup> 長谷川 剛<sup>††</sup> 村田 正幸<sup>†</sup>

<sup>†</sup> 大阪大学大学院情報科学研究科  
〒 560-0871 大阪府吹田市山田丘 1-5  
<sup>††</sup> 大阪大学サイバーメディアセンター  
〒 560-0043 大阪府豊中市待兼山町 1-32

E-mail: †{mlt-cao,murata}@ist.osaka-u.ac.jp, ††hasegawa@cmc.osaka-u.ac.jp

あらまし 高速 (1 ギガビットまたはそれ以上) ネットワークにおいて、パケットペアやパケットトレインなど、パケット間隔ベースの計測手法は、下記の二つの問題を持つ。まず、高速ネットワークでの計測は非常に短いパケット間隔が必要となるが、短い間隔でパケットを送信することは、大きな CPU 負荷を必要とする。次に、高速ネットワーク対応のネットワークインタフェースのほとんどが割り込み削減機構 (IC、Interrupt Coalescence) を採用していることである。IC はパケットの到着間隔を変えたり、パケットのバースト転送を生成したりするため、パケット間隔ベースの計測が不正確となる。そこで本稿では上記の二つの問題を解決した ICIM (Interrupt Coalescence-aware Inline Measurement) というインライン計測手法を提案する。インライン計測とは TCP の転送中のデータパケットを用いた計測手法である。提案手法は IC によって発生するパケットのバースト転送を逆に利用し、パケット送信間隔を調整することなく高い利用可能帯域を計測することができる。シミュレーション結果から、ICIM が IC がある環境で数 Gbps の利用可能帯域でも計測可能であることがわかった。さらに、既存のパケットストリームを用いた手法に比べて、計測に使用するパケットが 1/100 程度になることもわかった。また ICIM を導入した TCP が従来の TCP と同じデータ転送性能を持ちながら、数 RTT 程度という短い間隔で計測結果を導出することもわかった。

キーワード end-to-end 計測、利用可能帯域、インライン計測、パケットペア、パケットトレイン、パケットバースト

## An Inline Network Measurement Mechanism for High-Speed Networks

Cao LE THANH MAN<sup>†</sup>, Go HASEGAWA<sup>††</sup>, and Masayuki MURATA<sup>†</sup>

<sup>†</sup> Graduate School of Information Science and Technology, Osaka University  
1-3, Yamadagaoka, Suita, Osaka 560-0871, Japan  
<sup>††</sup> Cybermedia Center, Osaka University  
1-32, Machikaneyama, Toyonaka, Osaka 560-0043, Japan

E-mail: †{mlt-cao,murata}@ist.osaka-u.ac.jp, ††hasegawa@cmc.osaka-u.ac.jp

**Abstract** In high-speed networks, such as 1-Gbps or higher networks, bandwidth measurement algorithms that utilize packet transmission/arrival intervals, such as packet trains and packet pairs, have a number of problems. First, network measurement for large bandwidth requires short packet transmission intervals, which causes a heavy load on the CPU. Second, network interface cards for high-speed networks usually employ Interrupt Coalescence (IC), which rearranges the arrival intervals of packets and causes bursty transmission of packets. In the present study, we introduce ICIM (Interrupt Coalescence-aware inline measurement), a new bandwidth measurement approach that overcomes these two problems. ICIM utilizes the data packets of an active TCP connection for the measurement. In order to determine the available bandwidth, rather than adjusting the packet transmission intervals, the TCP sender instead adjusts the number of packets involved in a burst and checks whether the corresponding ACK packets also form a burst. Simulation results show that ICIM can measure the bandwidth as high as some Gbps while requiring a number of data packets that is only 1/100 of that of the existing stream-based algorithm. TCP that is deploying ICIM can yield measurement results in intervals as short as some RTTs, while maintaining the properties of the original TCP.

**Key words** end-to-end measurement, available bandwidth, inline measurement, packet pair, packet stream, packet burst

## 1. Introduction

Active measurement of the available bandwidth of an end-to-end network path has been vigorously investigated [1-3]. Compared with passive measurement, active measurement can deliver faster and more accurate results because the network can be investigated in detail using probe traffic. However, the sending of probe traffic is a drawback of active measurement. According to [3], Pathload [1] generated between 2.5 to 10 MB of probe traffic per measurement. Newer tools have succeeded in reducing the amount of probe traffic. The average per-measurement probe traffic generated by IGI [2] is 130 KB and that generated by Spruce [3] is 300 KB. Although a few KB of probe traffic for a single measurement is a negligible load on the network, for routing in overlay networks, or adaptive control in transmission protocols, these measurements may be repeated continuously and simultaneously from numerous network nodes and end hosts. In such cases, probe traffic of a few KB per measurement will generate a large amount of traffic that may interfere with data transmission in the network, as well as degrading the measurement itself.

Previously, we proposed an active measurement method that overcomes the problem mentioned above [4]. We proposed the concept of *inline measurement*, that is, the idea of “plugging” the active measurement mechanism into an active TCP connection. This method has the advantage of requiring no extra traffic to be sent on the network, and provides fast and accurate measurement. We refer to RenoTCP employing this mechanism as ImTCP (Inline measurement TCP). When the sender transmits data packets, ImTCP adjusts the transmission rate of some packets, and considering arrival intervals of the corresponding ACK packets, the ImTCP sender estimates the available bandwidth. ImTCP utilizes a measurement algorithm similar to that of Pathload [1]. That is, the arrival intervals of packets that are sent back-to-back at a specified rate are used to estimate the available bandwidth. ImTCP delivers measurement results repeatedly in short intervals, such as a few RTTs, and the number of packets involved in each measurement is far fewer than that for Pathload.

In the present study, we focus on a new challenge regarding active measurement. Specifically, we investigate the bandwidth measurement of 1-Gbps or faster network paths, which are becoming increasingly popular. In such high-speed networks, ImTCP, Pathload and other active measurement tools based on packet spacing [2, 3] must overcome the following problems. First, measurement in fast networks requires short transmission intervals of the probe packets (for example, 0.12 ms for a 1-Gbps link). However, regulating such short intervals causes a heavy load on the CPU. Second, network cards for high-speed networks usually employ Interrupt Coalescence (IC) [5, 6], which rearranges the arrival intervals of packets and causing bursty transmission, so that the algorithms utilizing the packet arrival intervals do not work properly.

We introduce a new inline measurement mechanism that works well in high-speed networks. We call this ICIM (Interrupt Coalescence-aware Inline Measurement). Unlike other active measurement tools, ICIM adjusts the number of packets that are transmitted in a burst caused by IC and estimates the available bandwidth by observing the number of packets in the burst as it passes through the network, rather than by observing the inter-intervals of the packets. ICIM does not set the sending interval of the packets, so the overhead for packet spacing at the sender is eliminated. The measurement results show that TCP with ICIM can transmit data with the same performance as Reno TCP and can measure the available bandwidth of high-speed networks.

The remainder of this paper is organized as follows. In Section 2, we discuss problems of measurement in high-speed networks and

look at a number of related studies. In Section 3, we introduce ICIM and explain how to realize it in Reno TCP. In Section 4, we evaluate the performance of Reno TCP that is utilizing ICIM. Finally, in Section 5, we present concluding remarks and discuss future projects.

## 2. Available bandwidth measurement in high-speed networks

In this section, we discuss some of the difficulties encountered by existing active available bandwidth measurement tools, including ImTCP, in high-speed networks (1 Gbps or higher). We assume that the machines that run the measurement tools are general purpose machines, for example, a x86-based CPU machine with a normal OS, such as 4.4 BSD LINUX. The problems mentioned here may not occur in high-performance machines that are designed especially for measurement.

### 2.1 Limitation of packet pacing in general-purpose machines

In current active measurement tools, probe packets must be sent at a rate higher than the available bandwidth of the network path, otherwise the packet space will not be expanded and the tools will not be able to determine the available bandwidth. When the available bandwidth can reach 1 Gbps or higher, the transmission intervals of the probe packets must be 0.012 ms (for measuring 1-Gbps bandwidth) or smaller. As we discuss below, for a general-purpose machine, sending packets in such small intervals causes high CPU overhead.

For pacing packets, there are two approaches. The first is to continuously check the hardware clock (for example, using `gettimeofday()` in UNIX systems) and send the packets when the clock reaches a determined timing. In a Linux system with an x86-based CPU, one access of the hardware clock requires approximately 1.9  $\mu s$  (in the FreeBSD system, one access requires 9  $\mu s$ ) [7]. The `write()` system call requires an average of 2  $\mu s$  (in the case of a Pentium III CPU). Therefore, a Linux system can only send packets in intervals greater than  $2 + 1.9 = 3.9 \mu s$ . This means that, the system can measure the bandwidth up to 3 Gbps (for the case in which the probe packet size is 1,500 Bytes). However, in order to send packets at 3 Gbps, the CPU has to spend all of the time checking the hardware clock overhead. If the measurement is repeated continuously, then the CPU will not be able to process tasks from other applications. The system performance then will be deteriorated. Thus, checking the hardware clock to send packets in a high-speed network is not a good approach.

The second approach is to register the packet sending program to an Interrupt Service Routine (ISR) of the hardware clock interrupt. In a general-purpose UNIX OS, the ISR `hardclock()` is provided for this purpose. In 4.4BSD OS and LINUX, the `hardclock()` system call is called by the interrupt of hardware clock every 0.01 s. However, with this low interrupt frequency, the program called by `hardclock()` can only send packets at the rate of 1.2 Mbps (assuming that the packet size is 1,500 Bytes). To obtain a higher interrupt frequency, a new interrupt schedule of the hardware clock can be implemented. However, one hardware interrupt (in 4.4 BSD OS) normally requires more than 1  $\mu s$  [8]. If the packet transmission rate is 1 Gbps, then the sending interval is 12  $\mu s$ . This means that, in this case, the overhead of the hardware interrupt is as high as 1/12 of the total working time of the CPU. In addition, a new interrupt schedule for the hardware clock requires many changes in the OS.

### 2.2 Effects of Interrupt Coalescence

Another reason that the task of measurement in high-speed networks difficult is IC, which is deployed in most high-bandwidth Network Interface Cards (NICs). IC is a technique in which NICs

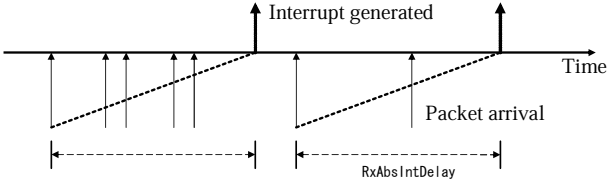


Fig. 1 Receive Absolute Timer

group multiple packets that arrive in a short time interval and pass them to the OS in a single interrupt. IC reduces the CPU overhead when the arrival intervals of packets become small. Because the inter-arrival intervals of the packets observed by the kernel are changed, IC has an enormous impact on bandwidth measurement tools, in which the arrival intervals of packets are utilized for bandwidth estimation.

There are a number of types of timer setting in IC. For example, Intel Gigabit Ethernet Controllers [5] contains the following mechanisms for IC:

- Absolute timer: The absolute timer delays the assertion of an interrupt to allow the controller to collect additional interrupt events before delivering them to software.
- Packet timer: The packet timers are inactivity timers, triggering interrupts when the link has been idle for an appropriately long interval.
- Master timer for throttling all interrupt sources: An interrupt throttling mechanism is used to set an upper bound for the interrupt rate.

The absolute timer is the default setting for Intel Gigabit Ethernet Controllers. Only users with root privileges can change the IC settings in NIC drivers. We investigate the absolute timers in greater detail. There are two absolute timers. One is for transmit interrupts, and the other is for receive interrupts. Because transmit interrupts only inform the kernel as to the completion of packet sending, delays in transmit interrupts do not affect the real transmission intervals of the packets. In contrast, delays in receive interrupts change the intervals of all receiving packets observed by the kernel. As shown in Figure 1, the receive absolute timer starts to count down upon receipt of the first packet. Subsequent packets do not alter the countdown. Once the timer reaches zero, the controllers generate an interrupt to pass all of the packets to the OS in a bursty manner. The length of the timer is decided by the parameter  $RxAbsIntDelay$ , which is defaulted to 0.1312 ms in Intel Gigabit Ethernet Controllers [9]. Thus, all packets that have time intervals smaller than  $RxAbsIntDelay$  will belong either to the same burst, in which case the time interval between the packets becomes zero, or to two successive bursts, in which case the time interval becomes  $RxAbsIntDelay$  or larger. Therefore, the software cannot detect packet intervals smaller than  $RxAbsIntDelay$ . With the default value of 0.1312 ms for  $RxAbsIntDelay$ , the software cannot perceive transmission rates larger than 100 Mbps (if the packet size is 1,500 Bytes).

There are some studies that have discussed measuring bandwidth using the existing IC. For example, one study [7] suggests that in order to obtain the real arrival intervals of packets, the onboard timestamp of some network cards (for example SysKconnect GigE NIC [6]) should be used. However, the same study also concludes that this solution is not useful for general-purpose network measurement tools, because very few NICs have an onboard timer. Furthermore, using an onboard NIC timer requires modification of the device driver. This prevents the tool from being easy to run on numerous systems.

Another study [10] reports that since the last packet in a burst formed by IC has the smallest delay in the NIC buffer, the intervals

of the last packets in the bursts can be used for estimation of the available bandwidth, according to the Pathload [1] algorithm. However, because only a small part of stream is used for the measurement, the stream must be very long. This is not suitable in inline measurement, because making long measurement streams in TCP badly affects the TCP transmission performance.

### 3. Interrupt Coalescence-aware Inline Measurement (ICIM)

#### 3.1 Effects of Interrupt Coalescence on TCP

The behavior of TCP when the network cards enable IC has been investigated in previous studies [8, 10], and IC has been shown to be detrimental to TCP self-clocking. IC causes the ACK packets to arrive at the sender in bursts, and this bursty arrival in turn causes bursty transmission of data packets and, subsequently, bursty transmission of ACK packets from the TCP receiver. According to one study [10], with IC, 65% of ACKs arrive with intervals of less than  $1 \mu s$ , because they are delivered to the kernel with a single interrupt. Meanwhile, without IC, almost no ACK packets arrive with small intervals.

In the present study, we propose an algorithm that can exploit the burst of data packets in TCP under the effects of IC to measure the available bandwidth. The TCP sender adjusts the number of packets involved in a burst and checks whether the corresponding ACK packets also form a burst to investigate the available bandwidth. ICIM can be employed into any version of TCP. Using previously reported results [10], ICIM first checks to see if the network card has IC enabled. If the IC is enabled, ICIM continues measurement based on the bursty transmission of TCP.

#### 3.2 Packet burst-based available-bandwidth measurement algorithm

Because the absolute timer (described in Section 2) is the default setting of the Intel(R) PRO/1000 Adapter [9], we assume that the NIC uses the absolute timer when receiving packets. The measurement algorithm using bursts of packets is described below.

As shown in Figure 2, we consider the situation in which two bursts of packets are sent at the interval  $S$ . The number of packets in the first burst (Burst 1) is  $N$ . Assume that  $C$  is the capacity of the bottleneck link.  $C_{Cross}$  is the average transmission rate of cross traffic over the bottleneck link when the two bursts pass the link, and  $P$  is the packet size. Then, the amount of traffic that enters the bottleneck link during the period from the point at which the first packet of Burst 1 reaches the link until the point at which the first packet of Burst 2 reaches the link will be the sum of the packets in Burst 1 and the cross traffic packets arriving in  $S$ , i.e.,  $C_{Cross} \cdot S + N \cdot P$ . If the amount is larger than the transfer ability of the link during this period, considered to be  $C \cdot S$ , then Burst 2 will go to the buffer of the link. This results in a tendency for the interval between the two bursts to increase after leaving the bottleneck link.

We can write that the burst interval will be increased if

$$C_{Cross} \cdot S + N \cdot P > C \cdot S \quad (1)$$

or,

$$\frac{N \cdot P}{S} > C - C_{Cross}$$

Note that  $C - C_{Cross}$  is the available bandwidth ( $A$ ) of the bottleneck link. Therefore, Eq. (1) becomes

$$\frac{N \cdot P}{S} > A$$

Since we assume that the absolute timer is used,  $S$  is always larger than  $RxAbsIntDelay$ . Therefore, at the NIC of the TCP receiver,

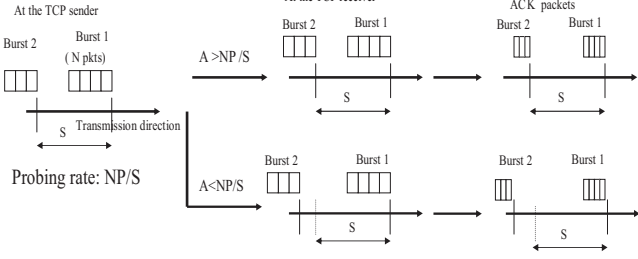


Fig. 2 Packet burst-based available-bandwidth measurement principle

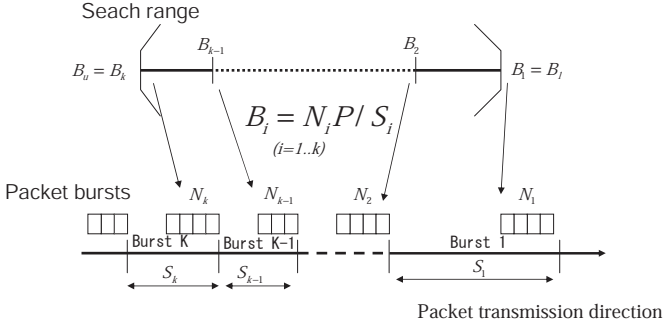


Fig. 3 Probing a search range in ICIM

since the arrival interval of the two bursts are larger or equal to  $S$ , the two bursts are passed to the kernel in two different interrupts. The TCP receiver then sends the ACK of the two bursts in the same intervals to the sender TCP. Thus, by checking the arrival intervals of the corresponding ACK packets of the two bursts, the TCP sender can determine if  $A > NP/S$ . By sending numerous bursts with various values of  $NP/S$  (by changing  $N$ ), we can search for the value of the available bandwidth  $A$ . This is the measurement principle of the proposed inline measurement mechanism.

### 3.3 ICIM

ICM inherits the concept of the *search range* from the measurement algorithm in ImTCP [4]. This is the idea of limiting the bandwidth measurement range using statistical information from previous measurement results rather than searching from 0 bps to the upper limit of the physical bandwidth for every measurement. By limiting the measurement range, we can keep the number of probe packets small.

At first, we explain how to search for the available bandwidth in a determined search range and then we present an overview of the measurement algorithm.

Assume that the search range for a measurement is  $(B_l, B_u)$ . The algorithm then check  $k$  values in the range to determine which is nearest to the real available bandwidth. We use  $k = 4$  in the following simulations. The  $k$  points are:

$$B_i = B_l + \frac{B_u - B_l}{k - 1}(i - 1) \quad (i = 1, \dots, k)$$

The TCP sender then sends  $k$  consequence bursts and the number of packets are adjusted so that the probe rate of Burst  $i$  is  $B_i$ :

$$\frac{N_i \cdot P}{S_i} = B_i \quad (2)$$

We illustrate the setting in Figure 3.

Realization of Eq. (2) requires the following:

- The value of  $S_i$  must be estimated at the timing of the transmission of Burst  $i$ . In fact,  $S_i$  is unknown until Burst  $i + 1$  is transmitted. But we need the value at the timing of the transmission of Burst  $i$  in order to guarantee Eq. (3). We therefore

estimate the value of  $S_i$  by assuming that the amount of data in Burst  $i$  is proportional to the length of the interval as follow:

$$S_i = \frac{N_i \cdot P}{T} \quad (3)$$

where  $T$  is the average throughput of TCP.

- In case the number of packets in Burst  $i$  is smaller than  $N_i$ , additional packets must be added to the burst so that the packet number becomes  $N_i$ . ICIM utilizes a buffer located at the bottom of the TCP layer in order to store the packets temporarily before sending them to the IP layer, in the manner of ImTCP. ICIM stores all of the packets of the burst that preceded Burst 1 in the buffer. Packets are added to Burst  $i$  ( $i = 1..k$ ) when necessary in order to maintain the desired number of packets ( $N_i$ ) in these bursts.

ICIM sends  $k$  bursts and checks the corresponding ACK of the bursts. If from burst number  $j$ ,  $j = 1..k$ , the arrival interval of the bursts becomes larger, then  $B_j$  is considered to be the value of the available bandwidth in that measurement.

The measurement algorithm of ICIM is as follows:

- (1) Check whether IC is enabled.  
ICIM first checks whether IC is enabled for the network card. For the reasons explained in Subsection 3.1, ICIM checks the arrival intervals of the ACK packets. If more than 50% of the intervals are less than  $1 \mu s$ , then ICIM decides that IC is enabled. If the IC is enabled, then ICIM continues the measurement. Otherwise, the measurement algorithm introduced in ImTCP is used.
- (2) Set the initial search range  
We set the initial search range as  $(T, 2 \cdot T)$  where  $T$  is the throughput of TCP.
- (3) Wait until the window size (*cwnd*) is larger than  $C_{min}$  (large enough to create bursts for measurement). We use  $C_{min} = 50$  in the following simulations. Data packets are then sent in order to search the available bandwidth in the decided search range, as described above.
- (4) Add the new measurement result to the database.  
Calculate the new search range  $(B'_l, B'_u)$  from the database of the measurement results. The search range is calculated as follow:

$$B'_l = R - \max\left(1.96 \frac{V}{\sqrt{q}}, \frac{R}{10}\right)$$

$$B'_u = R + \max\left(1.96 \frac{V}{\sqrt{q}}, \frac{R}{10}\right)$$

where  $R$  is the latest measurement result.  $V$  is the variance of stored values of the available bandwidth and  $q$  is the number of stored values.  $A/10$  is a value that ensures that the search range does not become too small. Moreover, when measurement result in Step 3 falls to  $B_l$  ( $B_u$ ), it is possible to consider that the network has changed greatly so that the real value of the available bandwidth is lower (higher) than the search range. In this case, we discard the accumulated measurement results because they become unreliable as statistic data and enlarge the search range  $(B_l, B_u)$  twice towards the lower (higher) direction to create  $(B'_l, B'_u)$ .

- (5) Wait for  $Q$  seconds then return to Step 2 and start the next measurement. During the waiting time  $Q$ , TCP transmits packets in the normal manner. The waiting time is needed for the TCP transmission to return to the normal state after the packets store-and-forward process at Step 2.

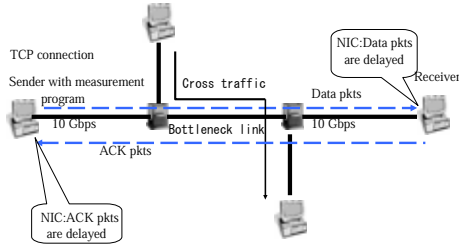


Fig. 4 Simulation topology

## 4. Simulation experiment

### 4.1 Measurement results

We show the measurement results for ICIM through ns-2 simulations. We deploy ICIM into Reno TCP, the most popular version of TCP, and use the topology shown in Figure 4 for the simulation. The sender and receiver of TCP are connected through 10 Gbps access links and a bottleneck link. The NICs of both the sender and receiver host employ IC with the absolute timer. The cross traffic on the bottleneck link is made up of UDP flows, in which various packet size are used according to the monitored results in the Internet reported in [11]. The Capacity of the bottleneck link is 5 Gbps and the available bandwidth (A-bw) is 2 Gbps (from 0 to 15 sec), 3 Gbps (from 15 to 35 sec) and 4 Gbps (from 35 to 50 sec).

Figures 5(a) and 5(b) show the measurement results for ICIM when the interval between two measurements  $Q$  are set to 1 RTT and 2 RTTs, respectively. The curved line "Average" shows the exponential moving average of the measurement results. Also shown are the search ranges for each measurements. We can see that the search ranges, in most of the cases, successfully cover the correct value of the A-bw. Therefore, ICIM can detect fast that value, even in such a high-speed network. When  $Q = 1$ , in every RTTs TCP sender stores and sends packet bursts for the purpose of measurement. These bursts makes the throughput of TCP slightly oscillate so that the estimation for the burst interval in Eq. (3) becomes wrong. Therefore, the probing rate of each Burst  $i$  may not be exactly  $B_i$  (in Step 3 of Section 3.3). This leads to a large dispersion of the measurement results in Figure 5(a). When  $Q = 2$ , TCP sender creates fewer packet bursts so the measurement results are nearer to the correct value of A-bw as shown in Figure 5(b). However, the measurement frequency (16.7 results/second) becomes only a half of that when  $Q = 1$  (34.2 results/second).

### 4.2 Comparison with IC-aware Pathload

For the comparison between ICIM and Pathload, the TCP sender and receivers are next replaced by the sender and receiver of Pathload. We use the version of Pathload that can detect and filter the effects of IC [10]. To make the measurement of Pathload faster, we set the starting probing rate to 200 Mbps (instead of 1 Mbps in default).  $\omega$  and  $\chi$  are set to 200 Mbps and 150 Mbps, respectively. We set the size of probing packets to 1500 bytes.

The measurement results of Pathload when the number of packet in a stream  $K$  is set to 160 are shown in Figure 6(a). Because the default value of  $RxAbsIntDelay$  used in NIC is 0.000132 (s) and the packet size is 1500 bytes, the number of packets in a burst, on the average, is 22 when the A-bw is 2 Gbps, 33 when A-bw is 3 Gbps and 44 when A-bw is 4 Gbps. Therefore, when  $K=160$ , there are about 9 bursts in each stream when A-bw is 2 Gbps. This means that Pathload has about 9 packets (the last one in the bursts) for the measurement. The increasing trend in the stream in this case can be determined well so Pathload can deliver good measurement results. However, when A-bw becomes 3 Gbps or larger, the number of bursts becomes about 6 or fewer. Pathload then does not have

Tab. 1 Comparison in number of packets required for a measurement

$A - bw$	ICIM	IC-aware Pathload	Ratio ICIM:Pathload
2 Gbps	110	$200 \cdot 12 \cdot 8 = 19\ 200$	0.006
3 Gbps	130	$200 \cdot 12 \cdot 9 = 21\ 600$	0.006
4 Gbps	154	$200 \cdot 12 \cdot 10 = 24\ 000$	0.006

Tab. 2 Throughput (Mbps) of Reno TCP using ICIM: Normal Reno TCP (ratio)

#connections	$Q=1$ RTT	$Q=2$ RTTs
4	466.4 : 490.6 (0.95:1)	483.7 : 475.6 (1.01:1)
8	451.1 : 544.4 (0.82:1)	505.1 : 490.5 (1.02:1)
12	418.7 : 577.7(0.72:1)	503.5 : 493.2 (1.02:1)

enough packets to detect well the increasing trend in the stream. Therefore, as shown in Figure 6(a), Pathload fails to deliver good measurement results when the bandwidth is equal to or larger than 3 Gbps.

Figure 6(b) shows the measurement results of Pathload when  $K$  is set to 200. In this case, Pathload have sufficient number of packets for detecting the increasing trend of streams therefore the measurement results are correct. However, as Pathload searches for the A-bw from a low value so it takes long time for yielding one result. The measurement frequency is only 0.28 results/second, 60 times smaller than that of ICIM (with  $Q=2$  RTTs).

Figure 6(b) shows that, if the A-bw changes during a measurement, Pathload may not detects the change well. At 15 second, the a-bw changes from 2 Gbps to 3 Gbps while Pathload is probing a rate smaller than 2 Gbps. When the probing rate reaches 2 Gbps, the A-bw is already changed, therefore Pathload can successfully detect the value 3 Gbps. However, at 35 second, the probing rate of the on going measurement reaches 3 Gbps before the change of A-bw from 3 to 4 Gbps so Pathload assumes that the A-bw is smaller than or equal to 3 Gbps. Therefore, Pathload delivers a value around 3 Gbps at the end of that measurement, that is far from the value of A-bw at this timing.

Table 1 compares the number of packets used in a measurement of ICIM and Pathload. ICIM sends four bursts of packets for each measurement. The average number of packets totally in four bursts are shown in the second column of the table. On the other hand, Pathload probes 8, 9, 10 times for one measurement results when A-bw is 2, 3, 4 Gbps, respectively. Each probe requires 12 streams, of which the number of packets is 200. We can see that the number of packets ICIM used is less than one percent of that of Pathload.

Figures 5 and 6 show that the measurement results of ICIM have a larger dispersion in comparison with Pathload. That is because, base on the nature of algorithm, ICIM cannot increase the length of each measurement burst to obtain high accuracy as Pathload does. Instead, the accuracy can be improved by taking the exponential moving average in suitable intervals.

### 4.3 TCP compatibility

We next examine the data transmission performance of Reno TCP when employing ICIM. We perform a simulation where a number of TCP connections using ICIM conflict with the same number of TCP connections, which is not using ICIM, through a 1 Gbps bottleneck link. All the connections have the same RTT (0.018 s) and the same access link's bandwidth (10 Gbps). The number of connections is set to 4, 8 and 12. For each value of connection numbers, the simulation is repeated 10 times and the throughput of the TCP connections that have and do not have ICIM (and the ratio of them), are calculated and compared.

Table 2 shows the results when  $Q$  of ICIM is set to 1 RTT and

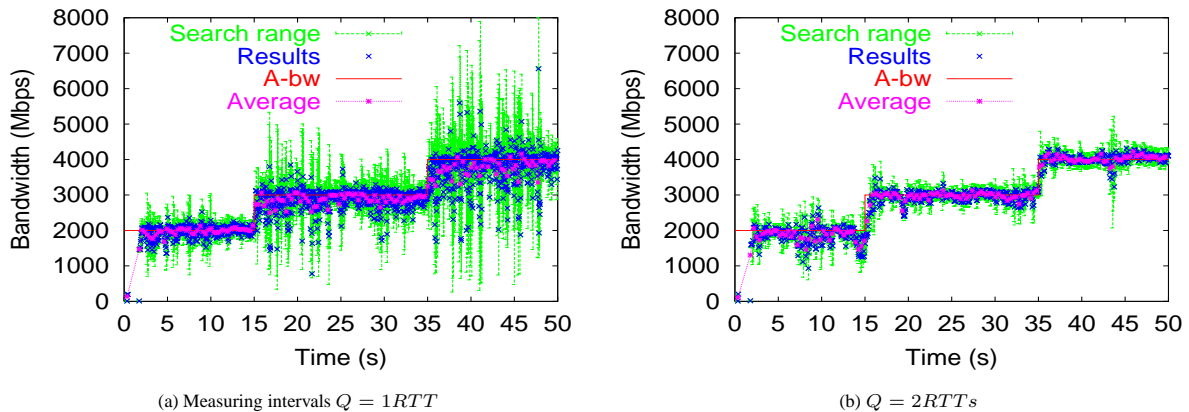


Fig. 5 Measurement results for ICIM

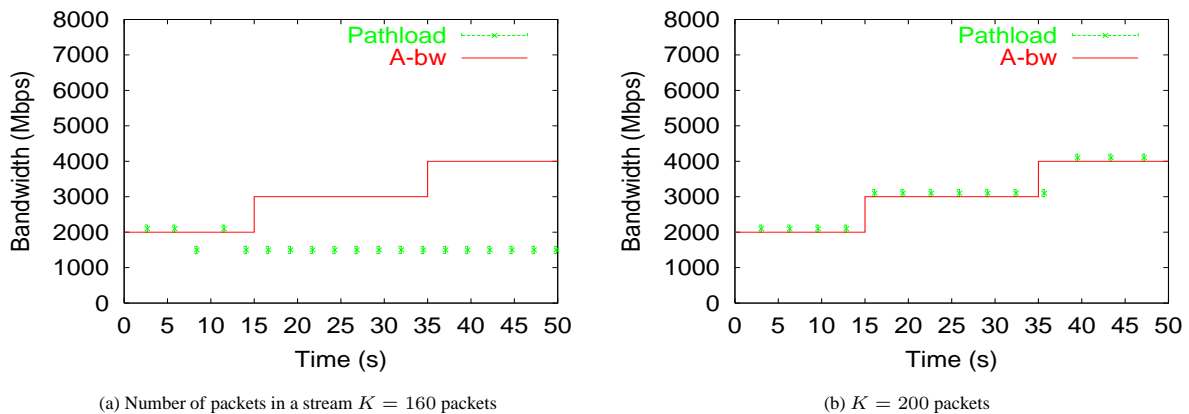


Fig. 6 Measurement results for IC-aware Pathload

2 RTTs. In case ICIM performs measurement in every RTT, the TCP archives lower throughput than TCP that do not perform ICIM when conflicting because ICIM has to delay many data packets for the measurement in this case. As shown in Table 2, the ratio of throughput between ImTCP HighSpeed over RenoTCP is always smaller than 1. When the number of connections increases, the ratio is lower because the conflict between TCP connections are more intensive. If ICIM takes a lower measurement frequency, for example, when  $Q = 2$  RTTs, then the TCP connections performing ICIM can obtain the same throughput with the normal Reno TCP, as shown in the third column of the table.

## 5. Conclusion and future studies

In the present paper, we introduced ICIM, a new method that can measure the available bandwidth in a 1-Gbps or higher network path. The proposed measurement algorithm does not require regulation of packet transmission intervals and works well with Interrupt Coalescence. Simulation experiments showed that the proposed measurement algorithm works well with no degradation of TCP data transmission speed.

At present, we are evaluating the performance of ICIM in a real network environment. In addition, we are investigating the measurement mechanism for the capacity of high-speed networks that can be implemented in ICIM with the least change.

### References

- [1] M. Jain and C. Dovrolis, "End-to-end available bandwidth: Measurement methodology, dynamics, and relation with TCP throughput," in *Proceedings of ACM SIGCOMM 2002*, Aug. 2002.
- [2] N. Hu and P. Steenkiste, "Evaluation and characterization of available bandwidth probing techniques," *IEEE Journal on Selected Areas in*

- Communications*, vol. 21, Aug. 2003.
- [3] J. Strauss, D. Katabi and F. Kaashoek, "A measurement study of available bandwidth estimation tools," in *Proceedings of Internet Measurement Conference 2003*, Oct. 2003.
- [4] Cao Le Thanh Man, Go Hasegawa and Masayuki Murata, "Available bandwidth measurement via TCP connection," in *Proceedings of the 2nd Workshop on End-to-End Monitoring Techniques and Services E2EMON*, Oct. 2004.
- [5] Intel, "Interrupt Moderation Using Intel Gigabit Ethernet Controllers," available at <http://www.intel.com/design/network/applnots/ap450.pdf> (2003).
- [6] Syskonnect, "SK-NET GE Gigabit Ethernet Server Adapter," available at [http://www.syskonnect.com/syskonnect/technology/SK-NET\\_GE.PDF](http://www.syskonnect.com/syskonnect/technology/SK-NET_GE.PDF) (2003).
- [7] G. Jin and B. Tierney, "System capability effect on algorithms for network bandwidth measurement," in *Proceedings of Internet Measurement Conference 2003*, Oct. 2003.
- [8] M. Zec, M. Mikuc and M. Zagar, "Estimating the impact of interrupt coalescing delays on steady state TCP," in *Proceedings of the 10th SoftCOM 2002 conference*, 2002.
- [9] Intel(R) PRO/1000 Adapter, "README file," available at [http://support.intel.co.jp/jp/support/network/adapter/1000/linux\\_readme.htm](http://support.intel.co.jp/jp/support/network/adapter/1000/linux_readme.htm).
- [10] R. Prasad, M. Jain and C. Dovrolis, "Effects of interrupt coalescence on network measurements," in *Proceedings of the 5th Passive and Active Measurement Workshop PAM 2004*, Apr. 2004.
- [11] "NLNLR web site," available at <http://moat.nlanr.net/Datacube/>.