

# On Detecting Congestion Signals from Relative One-Way Delay

Yoshito Tobe<sup>1</sup> Hiroto Aida<sup>2</sup> Yosuke Tamura<sup>1</sup> Hideyuki Tokuda<sup>1,2</sup>

<sup>1</sup> Graduate School of Media and Governance <sup>2</sup> Faculty of Environmental Information  
Keio University  
5322, Endo, Fujisawa-shi, Kanagawa, 252-8520 Japan  
{tobe,haru,tamura,hxt}@ht.sfc.keio.ac.jp

## Abstract

In this paper, we describe the experimental results of measuring per-packet *Relative One-way Trip Time* (ROTT) of a UDP stream at a receiver on the Internet. It is found that spikes with successive plots, which we call *spike-trains*, often appear on a time-ROTT graph. Also, congestion-related losses are found to be strongly correlated only to the spike-trains and the path status is effectively identified by such spike-trains. This paper is a preliminary work towards detecting congestion signals from ROTT.

**Keywords:** delay, measurement, congestion control, Internet

## 1 Introduction

Recent advances in computer technology, digital signal processing applications, and packet network capacity have accelerated the proliferation of real-time, packet-based continuous media (CM) communications. Furthermore, the growth of the Internet has fueled demand for these technologies and highlighted the difficulties of handling CM in a wide area network. Once a CM flow is injected into a best-effort Internet, a rate control that avoids congestion is essential. UDP flows that carry CM do not automatically control their rates, which may make both the network and themselves inefficient. To overcome this problem, TCP-friendly rate control algorithms [7, 9, 13, 8] have been introduced.

A problem in the rate control is the type of information whereby congestion is detected and determined. With the exception of some proposals [3, 1], TCP relies on a loss of acknowledgments (ACKs) indicating a packet loss to detect congestion. Unlike a TCP flow, a sender of a CM UDP flow does not usually receive per-packet ACKs. Instead, it is recommended that the CM UDP flow uses Real-time Transport Protocol (RTP) and its control protocol, RTCP, to control a rate based on information about losses in a certain period. However, reliance on information about losses alone may lead to an erroneous control. First, when the UDP flow competes with a TCP flow, the TCP flow reduces its rate more quickly, resulting in a few or even no packet losses in the UDP flow. In other words, reliance on the information about loss only works when a UDP receiver notifies the corresponding sender of the loss immediately and the sender reacts to the notification quickly. Second, as is well known, packets may be lost in the absence of congestion over a lossy wire-

less link [4].

To explore an alternative way to detecting congestion, we investigate a way to measure one-way delay from the sender to the receiver together with losses. Since an absolute value of delay is difficult to obtain due to a skew between the clocks at the sender and the receiver, we observe variation in the one-way delay. We call such a variation *Relative One-way Trip Time* (ROTT). Extensive experimental results show that a value of ROTT often increases with a spike accompanied by successive plots, which we call *spike-trains*, on a time - ROTT graph.

We then investigate the correlation between the ROTT and losses, and find that spike-trains are only related to congestion-related losses.

The remainder of this paper is organized as follows. In section 2, some measured results of ROTT over the Internet are shown. In section 3, the dependence between ROTT and packet losses is analyzed. In section 4, related works are described.

## 2 Path Status

### 2.1 Identification of the Variation in One-Way Delay

Prior to discussing how to determine the path status, we define how to deal with delay. Since absolute values of one-way delay are difficult to measure, we calculate the variation in delay, instead. First, we assume that a sender is transmitting rate-controlled data at a certain regular interval with a header including a timestamp and a sequence number. Let  $t_s(i)$  ( $i = 1, 2, 3, \dots$ ) and  $t_r(i)$  denote the time that  $i$ -th packet of the CM flow is sent and received, respectively. Then we define a *relative delay*  $\Delta t_{r,s}(i)$

as follows:

$$\Delta t_{rs}(i) \equiv t_r(i) - t_s(i) - (t_r(1) - t_s(1)).$$

Note that  $t_s(i)$  and  $t_r(i)$  are measured with the clock of the sender and the receiver, respectively. Hence  $\Delta t_{rs}(i)$  monotonically increases or decreases depending on the difference between the precise tick of these two clocks. Therefore,  $\Delta t_{rs}(i)$  is expressed as follows:

$$\Delta t_{rs}(i) = D_{min} + \Delta T_{rs}(i) + \alpha t(i) + \beta,$$

where  $t(i)$ ,  $\alpha$ , and  $\beta$  are the receiving time of  $i$ -th packet, the skew and the offset caused by the difference between the clocks. Here,  $D_{min}$  which is referred to as the base delay is the minimum value of true one-way trip time. However we cannot obtain the value of  $D_{min}$  unless the clocks of the sender and the receiver are completely synchronized. Instead,  $\Delta T_{rs}(i)$  which is referred to as *Relative One-way Trip Time* (ROTT) can be obtained if  $\alpha$  and  $\beta$  are estimated. We use ESRS [11] to calculate  $\alpha$  and  $\beta$ <sup>1</sup>.

## 2.2 Basic Observation of Delay and Loss

We examine the relationship between the loss and ROTT through experiments over the Internet (Figure 1). Communication paths between the following sites were used: **Site1 in Spain**, **Site2 in U.S.A.**, and **Site3 in Japan**<sup>2</sup>. Typical RTTs on the Site3-Site2 and Site1-Site3 paths are 140-300 ms and 700-1500 ms. All hosts are installed to FreeBSD2.2.8R and equipped with a Pentium counter for measurement. The timestamp contained in the header uses a value converted from the Pentium counter. As a result, it does not mean the absolute time. In this paper, a rank-1 rate denotes a rate which is obtained by averaging a rate over 1 s.

First we used the path between Site2 and Site3. We created one flow of TCP from Host A1 to B and another of Rate-probing UDP (RUDP) from Host A2 to B simultaneously. In RUDP, rate probing using TCP is intermittently inserted [10]. Both RUDP and TCP flows transmit 1440-byte packets<sup>3</sup>.

Figure 2 shows rank-1 rates of RUDP and TCP flows measured at Host B. We will examine the cause of frequent dropping in rank-1 rate of TCP by investigating ROTT. We then examine ROTT of the RUDP flow. In Figure 3, we set the point at which the first rate probing finishes to the beginning of calculating ROTT. As seen in the figure, several jumped sequences of plots appear in the delay. We

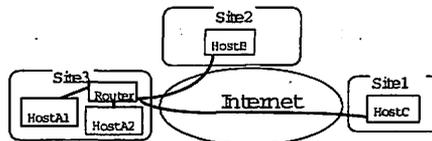


Figure 1: Topology over the Internet

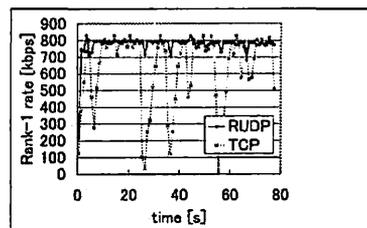


Figure 2: rank-1 rate (Site3 - Site2)

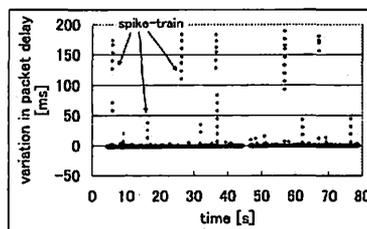


Figure 3: Spike-trains appears in ROTT (Site3 - Site2)

call such a sequence a “spike-train.” This is the same phenomenon as probe compression in [2], but it is more clearly seen in a time - ROTT graph. Comparing the trend of the delay with Figure 2, a dropping in the TCP rank-1 rate occurs when a spike-train appears. This suggests that packets are likely to be dropped when a spike-train appears.

Let us further see where packet losses occur. Figures 4 - 6 show enlarged parts of Figure 3. As can be seen in these figures, the delay does not always increase when a packet is lost. At the same time, a packet is not always lost when a spike-train appears. However, most packets are lost in spike-trains. We will examine the relationship further later.

We also conducted an experiment over the path between Site1 and Site3. One 1440-byte UDP data was transmitted at a regular interval. As seen in Figures 7 and 8, spike-trains are also observed in the Site1-Site3 path.

<sup>1</sup>The details of ESRS are described in [11].

<sup>2</sup>Site1, Site2, and Site3 are CESAT, Carnegie Mellon University, and Keio University, respectively.

<sup>3</sup>The size of TCP or UDP payload was set to 1440 bytes, which is a default TCP maximum segment size over Ethernet in FreeBSD 2.2.8R.

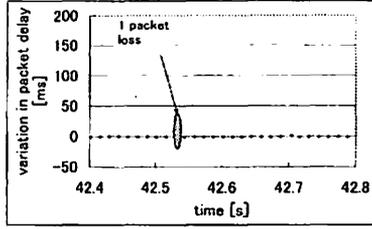


Figure 4: One packet loss without variation (Site3 - Site2)

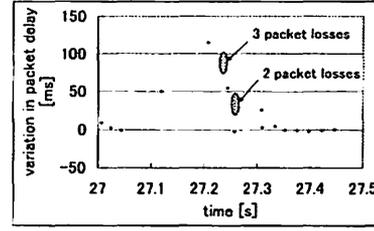


Figure 8: Packet losses in a spike-train (Site1 - Site3)

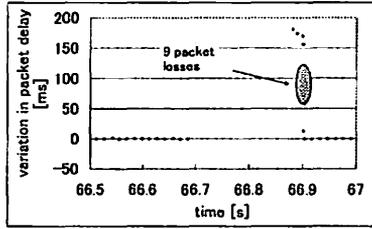


Figure 5: Packet losses in a spike-train (Site3 - Site2)

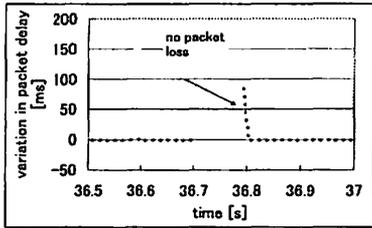


Figure 6: No packet loss in a spike-train (Site3 - Site2)

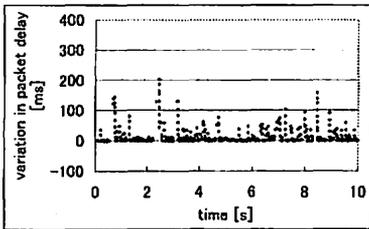


Figure 7: Spike-trains (Site1 - Site3)

### 3 Dependence between ROTT and Loss

Based on the preliminary observations above, we further investigate the dependence between ROTT and loss, and also examine the effect of reducing a rate.

Let us define the variables for analyzing the relationship among ROTT, loss, and rate as listed in Table 1. We created Flows 1 and 2 transmitting 1440-byte packets at a regular interval with different rates simultaneously on the path between Site2 and Site3. The transmission rate of Flow 2 is varied to six cases, 200 kbps, 400 kbps, 900 kbps, 1200 kbps, 1800 kbps, and 3000 kbps, whereas the transmission rate of Flow 1 is fixed at 600 kbps. Ten trials with 80-s duration were conducted for six cases and the variables in Table 1 were calculated.  $N_p$  is identified by the difference between the start and end of sequence numbers. To calculate  $N_s$ , an algorithm for detecting a spike-train shown in Figure 9 was used. The threshold values  $B_{spikestart}$  and  $B_{spikeend}$  were set to 20 ms and 5 ms, respectively.

Table 2 shows the results obtained. Since an actual transmitted rate is calculated from  $N_p$ , the ratio of rates for the two flows  $\gamma$  is calculated by  $N_p(\text{Flow 1}) / N_p(\text{Flow 2})$ . Each case was conducted in a different time. However, a common characteristics can be identified from the table. First,  $\rho_s$  is almost the same in the two flows. This suggests that a flow suffers from spike-trains in the identical ratio, irrespective of its transmission rate. Second,  $\rho_{s,l}$  is always near 1.0; most packet losses occur in the spike-train periods. There are some exceptions, but the ratio of them is negligibly small. Finally,  $\rho_{s,s}$  is approximately identical for the two flows. A possible explanation is that a flow with a larger rate may suffer from a larger number of losses but it can inject more packets in a spike-train period.

It was found that the maximum TCP-equivalent rate was always less than 800 kbps. Nevertheless, flows with a higher rate than 800 kbps, even with 3000 kbps, do not encounter packet losses at a higher rate. This indicates that a rate reduction scheme for

```

=====
On receipt of a packet with sequence number i:
  if ((state == not-in-spike-train) AND (ROTT(i) ≥ Bspikestart)) {
    state = in-spike-train;
    spike-first-sequence-number = i;
  }
  if ((state == in-spike-train) AND (ROTT(i) ≤ Bspikeend)) {
    state = not-in-spike-train;
    number-of-spike-packets = i - spike-first-sequence-number + 1;
  }
=====

```

Figure 9: Calculating the number of packets in a spike-train period

a RTP/UDP cannot rely on information about the loss ratio alone.

Table 1: Variables for analysis

$N_p$	number of packets expected to be received
$N_l$	number of lost packets
$N_s$	number of packets in spike-train periods
$N_{sl}$	number of lost packets in spike-train periods
$N_{bl}$	number of lost packets in blackouts
$\rho_{pl}$	$N_l/N_p$
$\rho_s$	$N_s/N_p$
$\rho_{sl}$	$N_{sl}/(N_l - N_{bl})$
$\rho_{ss}$	$N_{sl}/N_s$

Similarly, we measured values of the variables in Table 1 for one UDP flow of 64-byte packet on the Site1 - Site3 path. The transmission rate was set to 10 kbps because the maximum TCP rate was measured to be less than 30 kbps. Results of four trials with 80-s duration are shown in Table 3. In addition, snapshots of ROTT for trials 1 and 4 are shown in Figures 10 and 11, respectively.

In trials 1 and 3, blackout losses appear. In contrast, in trials 2 and 4, there is no such blackout but most packets are in spike-train periods. Hence it may be inappropriate to transmit a CM UDP data on such an ill-conditioned path. However, it is interesting to find that  $\rho_{ss}$  is almost the same in all cases in this experiments on the path.

The differences between the Site3-Site2 and Site1-Site3 paths are clearly observed in phase plots [2]. When a path is stable with a few spike-trains, sequences of plots both beside the vertical line and along the line

$$\Delta T_{r,s}(i+1) = \Delta T_{r,s}(i) - \Delta H(\text{constant}) \quad \dots(a)$$

are evident as shown in Figures 12 and 13. In contrast, when a path is unstable with many overlapped spike-trains, there are no such evident plot (Figure 14).

According to [2],  $\Delta H$  in Equ. (a) is related to an interval of transmission  $\delta$ ;  $\Delta H = \delta - P/\mu$ , where  $P$  and  $\mu$  are the length of packets and the ser-

vice rate at the bottleneck, respectively. Let  $P_p$ <sup>4</sup> and  $u$  denote the size of packet payloads and the transmission rate, respectively. Since  $\delta = P_p/u$ ,  $\mu = P/(P_p/u - \Delta H)$ . This is converted to payload-level bottleneck bandwidth  $\mu'$ ;  $\mu' = P_p/P \times \mu = P_p/(P_p/u - \Delta H)$ . Thus, a spike-train contains information about the bottleneck bandwidth. For instance, in a trial in Set-6 Site3-Site2 experiment,  $\Delta H$  of Flow 2 was measured to be approximately 3.0 ms. Since  $P_p = 1440$  byte, and  $u = 2674$  kbps,  $\mu$  is approximately 8.8 Mbps, which is much larger than the rate of Flow 2, 2675 kbps, and a TCP-equivalent rate, 800 kbps. This suggests that a rate control based only on the bottleneck bandwidth is difficult when the path is relatively stable.

Rather than an absolute value of  $\Delta H$ , the sign of  $\Delta H$  provides more useful information. It tells whether the transmission rate exceeds the bottleneck bandwidth. When  $u > P_p/(P/\mu + \Delta H)$ ,  $\Delta H > 0$ . In Figures 11 and 14, there are many plots where  $\Delta H > 0$ ; the transmission rate should have been set to be below 10 kbps. To observe the distribution of  $(\Delta T_{r,s}(i+1) - \Delta T_{r,s}(i))$ , its cumulative distribution functions (CDFs) are shown in Figures 15 - 17. We focus on the area where ROTT  $\geq 10$  ms because the behaviors only when ROTT increases are concerned with congestion. As can be seen in the graphs, most values of  $(\Delta T_{r,s}(i+1) - \Delta T_{r,s}(i))$  are distributed in the negative area for the Site3-Site2 path, and thus the CDF provides indication of stability of the transmission rate.

### 3.1 TCP Rate vs. ROTT Change

We also examine how the existence of spike-trains affect a TCP rank-1 rate. Again, Flows 1 and 2 are simultaneously run on the Site3-Site2 path except that Flow 2 is set to TCP. Let  $H$  and  $\eta_d$  denote the largest height of spike-trains in Flow 1 and the ratio of dropping in a rank-1 rate for Flow 2 in one-second period, respectively. The obtained results are shown in Figure 18. Interestingly, plots are grouped into two, which suggests that there are two major

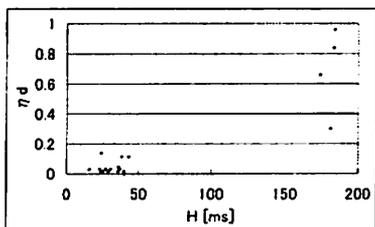
<sup>4</sup> $P = P_p + \text{size of UDP and IP headers.}$

Table 2: Measured statistics (Site3 - Site2)

Set		$N_p$	$N_l$	$N_s$	$N_{sl}$	$N_{bl}$	$\rho_{pl}$	$\rho_s$	$\rho_{sl}$	$\rho_{ss}$	$\gamma$
1	Flow 1 (600 kbps)	39137	22	267	20	0	$5.6 \times 10^{-4}$	$6.8 \times 10^{-3}$	0.91	0.075	0.36
	Flow 2 (200 kbps)	14272	1	73	0	0	$7.0 \times 10^{-5}$	$5.1 \times 10^{-3}$	-	-	
2	Flow 1 (600 kbps)	39167	19	209	17	0	$4.9 \times 10^{-4}$	$5.3 \times 10^{-3}$	0.89	0.081	0.72
	Flow 2 (400 kbps)	26021	20	169	17	0	$7.1 \times 10^{-4}$	$6.0 \times 10^{-3}$	0.85	0.10	
3	Flow 1 (600 kbps)	39165	51	239	49	0	$1.3 \times 10^{-3}$	$6.1 \times 10^{-3}$	0.96	0.21	1.6
	Flow 2 (900 kbps)	62425	107	317	106	0	$1.7 \times 10^{-3}$	$5.1 \times 10^{-3}$	0.99	0.33	
4	Flow 1 (600 kbps)	39061	56	348	53	0	$1.4 \times 10^{-3}$	$8.9 \times 10^{-3}$	0.95	0.15	2.1
	Flow 2 (1200 kbps)	83817	218	761	216	0	$2.6 \times 10^{-3}$	$9.0 \times 10^{-3}$	0.99	0.28	
5	Flow 1 (600 kbps)	39184	4	194	3	0	$1.0 \times 10^{-4}$	$5.0 \times 10^{-3}$	0.75	0.015	3.2
	Flow 2 (1800 kbps)	125012	10	512	8	0	$8.0 \times 10^{-6}$	$4.1 \times 10^{-3}$	0.80	0.016	
6	Flow 1 (600 kbps)	39235	25	245	25	0	$6.4 \times 10^{-4}$	$6.2 \times 10^{-3}$	1.0	0.10	4.7
	Flow 2 (3000 kbps)	165750	110	910	98	0	$5.9 \times 10^{-4}$	$4.9 \times 10^{-3}$	0.80	0.10	

Table 3: Measured results (Site1 - Site3)

Trial	$N_p$	$N_l$	$N_s$	$N_{sl}$	$N_{bl}$	$\rho_{pl}$	$\rho_s$	$\rho_{sl}$	$\rho_{ss}$
1	1706	409	562	53	313	0.24	0.33	0.55	0.094
2	1693	143	1258	111	0	0.084	0.74	0.78	0.088
3	1699	472	741	75	243	0.28	0.44	0.80	0.101
4	1698	155	1611	149	0	0.091	0.95	0.96	0.092

Figure 18:  $H$  vs.  $\eta_d$ 

bottlenecks in the Site3-Site2 path. It is found that the rank-1 rate is not substantially degraded when  $H$  is less than 50 ms.

## 4 Related Work

Kim *et al.* [6] proposed a scheme named LIMD/H. In LIMD/H, a sender maintains the history of losses and distinguishes between congestion-related and non-congestion-related packet losses. The distinction is reflected in change of a decrease factor of AIMD. In [1], a TCP receiver distinguishes between the two using inter-arrival times. These approaches aim at less than 1-s rate. In contrast, we aim at a control at 10-s rate that prevents a receiver from sending frequent feedback to a sender.

Additionally, there have been several efforts to introduce delay into congestion control. Jain [5] advocated congestion avoidance based on a change in RTT. This work provides a good starting point for a consideration of delay. In TCP Vegas [3], RTT is used not only to adjust timeout values, but also to estimate an expected throughput. The difference

between an actual throughput and an expected one provides a base for rate control. However issues of stability and coexistence among other TCP variants are still being studied. LDA [9] uses RTT for its control, but LDA only obtains samples of RTT by exchanging RTCP messages and the sampled RTT values do not contain dynamic behaviors of the delay. Bolot, in the literature [2], studied properties of RTT on the Internet with UDP probing packets at a regular interval. This work describes several significant findings: the relationship between the interval and two consecutive RTT values, and existence of probe compression. The probe compression is similar to a spike-train; our study does not use packets only for probing and focuses on ROTT instead of RTT. This work suggests that spike-trains do not appear when the regular interval is sufficiently large. We think that, unlike in probing, a CM UDP flow that attempts to obtain as much bandwidth as possible is likely to encounter a spike-train.

## 5 Conclusion

This paper has described the experimental results of measuring per-packet ROTT of a UDP stream. We have measured ROTT over two paths, Japan - U.S.A. and Japan - Spain, and find that a spike-train provides a congestion signal.

It is also found that reducing the transmission rate does not effectively result in decrease of packet loss ratio. This suggests that an algorithm that decreases the transmission rate based on packet loss ratio alone is not effective under some level of rate. Instead of packet losses, ROTT is observed to be effective for extracting congestion signals. For a future work, we plan to investigate the cause of spike-trains and also to establish a congestion control scheme based on measuring ROTT [12].

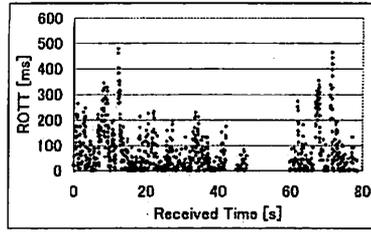


Figure 10: ROTT (Site1 - Site3, trial 1)

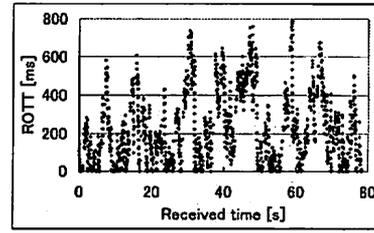


Figure 11: ROTT (Site1 - Site3, trial 4)

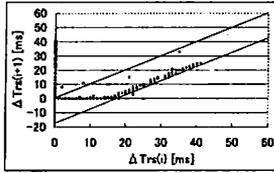


Figure 12: Phase plot of ROTT for case 5 (Site3-Site2:Flow 1)

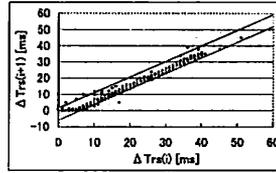


Figure 13: Phase plot of ROTT for case 5 (Site3-Site2:Flow 2)

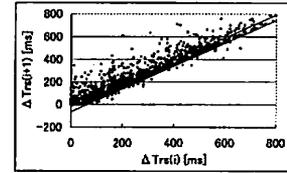


Figure 14: Phase plot of ROTT for trial 4 (Site1-Site3)

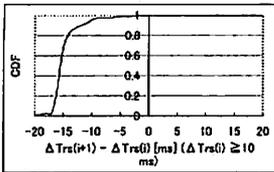


Figure 15: CDF of  $(\Delta T_{rs}(i+1) - \Delta T_{rs}(i))$  for case 5 (Site3-Site2:Flow 1)

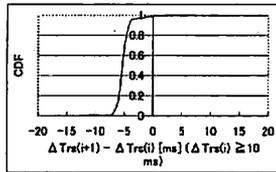


Figure 16: CDF of  $(\Delta T_{rs}(i+1) - \Delta T_{rs}(i))$  for case 5 (Site3-Site2:Flow 2)

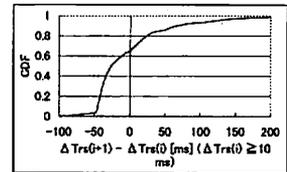


Figure 17: CDF of  $(\Delta T_{rs}(i+1) - \Delta T_{rs}(i))$  for trial 4 (Site1-Site3)

## References

- [1] S. Biaz and N. H. Vaidya, "Discriminating congestion losses from wireless losses using inter-arrival times at the receiver," *IEEE Symp. ASSET'99*, 1999.
- [2] J.-C. Bolot, "End-to-end packet delay and loss behavior in the Internet," *Proc. of ACM SIGCOMM'93*, pp. 289 - 298, Sept. 1993.
- [3] L. S. Brakmo, S. W. O'Malley, and L. L. Peterson, "TCP Vegas: new techniques for congestion detection and avoidance," *Proc. of ACM SIGCOMM'94*, pp. 24 - 35, Sept. 1994.
- [4] A. DeSimone, M. C. Chuah, and O. C. Yue, "Throughput performance of transport-layer protocols over wireless LANs," *Proc. of GLOBECOM'93*, Dec. 1993.
- [5] R. Jain, "A delay-based approach for congestion avoidance in interconnected heterogeneous computer networks," *ACM Computer Communication Review*, 19(5), pp. 56-71, Oct. 1989.
- [6] T. Kim, S. Lu, and V. Bharghavan, "Improving congestion control performance through loss differentiation," *Proc. of IEEE ICCCN'99*, Oct. 1999.
- [7] J. Padhye, J. Kurose, D. Towsley, and R. Koodli, "A model based TCP-friendly rate control protocol," *Proc. of the 9th Int. Workshop on Network and Operating Systems Support for Digital Audio and Video*, pp. 137-151, June 1999.
- [8] R. Rejaie, M. Handley, and D. Estrin, "An end-to-end rate-based congestion control mechanism for realtime streams in the Internet," *Proc. of INFOCOM'99*, Mar. 1999.
- [9] D. Sisalem and H. Schulzrinne, "The loss-delay based adjustment algorithm: a TCP-friendly adaptation scheme," *Proc. of the 8th Int. Workshop on Network and Operating Systems Support for Digital Audio and Video*, pp. 215-226, July 1998.
- [10] Y. Tobe, Y. Tamura, H. Nishino, and H. Tokuda, "Rate probing based adaptation at end hosts," *Proc. of IEEE Workshop on QoS Support for Real-Time Internet Applications*, pp. 58 - 65, June 1999.
- [11] Y. Tobe, H. Aida, Y. Tamura, and H. Tokuda, "Detection of change in one-way delay for analyzing the path status," *Proc. of the First Passive and Active Measurement Workshop*, ([http://pam2000.cs.waikato.ac.nz/final\\_program.htm](http://pam2000.cs.waikato.ac.nz/final_program.htm)), pp. 61 - 68, Apr. 2000.
- [12] Y. Tobe, Y. Tamura, A. Molano, S. Ghosh, and H. Tokuda, "Achieving moderate fairness for UDP flows by path-status classification," *Proc. of IEEE Local Computer Networks (LCN2000)*, Nov. 2000.
- [13] L. Vicisano, L. Rizzo, and J. Crowcroft, "TCP-like congestion control for layered multicast data transfer," *Proc. of INFOCOM'98*, Apr. 1998.