

WTCP: A Reliable Transport Protocol for Wireless Wide-Area Networks

Prasun Sinha Narayanan Venkitaraman Raghupathy Sivakumar Vaduvur Bharghavan

Coordinated Science Laboratory

University of Illinois at Urbana-Champaign

Email:{prasun, murali, sivakumr, bharghav}@timely.crhc.uiuc.edu

Abstract

Wireless wide-area networks (WWANs) are characterized by very low and variable bandwidths, very high and variable delays, significant non-congestion related loss, asymmetric up-link and downlink channels, and occasional blackouts. Additionally, the majority of the latency in a WWAN connection is incurred over the wireless link. Under such operating conditions, most contemporary wireless TCP algorithms do not perform very well.

In this paper, we present WTCP, a reliable transport protocol that is designed to operate efficiently and fairly over commercial WWAN networks such as CDPD. WTCP is rate-based, uses only end-to-end mechanisms, performs rate control at the receiver, and uses inter-packet delays as the primary metric for rate control.

We have implemented and evaluated WTCP over the CDPD network, and also simulated it in the ns-2 simulator. Our initial results indicate the WTCP can improve on the performance of comparable algorithms such as TCP-NewReno, TCP-Vegas, and Snoop-TCP by between 20% to 200% for typical operating conditions.

1 Introduction

Recent years have witnessed an explosive growth in the use of wireless wide-area networks (WWANs) such as CDPD, RAM and Ardis [1], with industry projections topping \$2.5 billion by the year 2002. Many large corporations are equipping their mobile workforce with laptops that have WWAN connectivity, thus enabling users to be connected anytime, anywhere. In the typical WWAN deployment scenario, a mobile user connects over the WWAN network to a dedicated proxy in the corporate backbone, which then acts as the service point for all the requests from the mobile user¹. Providing efficient and reliable connectivity between the proxy and the mobile host over commercial wide-area wireless networks is thus becoming a critical issue.

¹Most of the current deployment uses the proxy model for three reasons: (a) the connection quality of the WWAN network is too poor to sustain typical client-server applications, (b) the amount of data transferred to the mobile host must be filtered because of the orders of magnitude difference in bandwidth between the wired and WWAN connections, and (c) many portable computing devices have display and processing limitations that must be addressed by the proxy before sending the processed/filtered response to the mobile user.

Despite the enormous commercial interest, a typical connection over a WWAN network observes very low and highly variable throughput (between $O(100bps)$ – $O(10Kbps)$), very high and highly variable latency (between $O(400ms)$ – $O(4sec)$), bursty and random packet losses unrelated to congestion (1% - 10%), and occasional blackouts exceeding 10 seconds. Under such operating conditions, standard transport protocols such as TCP perform very poorly because of two main reasons: first, TCP assigns all packet losses to congestion and throttles down the transmission rate upon detecting a packet loss; and second, TCP sets its retransmission timeout (RTO) based on observed round trip times ($RTO = \overline{rtt} + 2 \times \Delta(rtt)$). In WWAN networks, neither mechanism is appropriate. Packet losses may happen either due to adverse channel conditions or congestion, and throttling down the transmission rate in the former case is the wrong thing to do. Furthermore, due to the very low bandwidth, transmission delays contribute significantly to the observed round trip time. Thus, round trip time computations are highly susceptible to fluctuation depending on the size of the TCP congestion window. As a consequence, RTO values computed by the sender tend to be very high - as large as 32 seconds for CDPD. A combination of large RTO for loss recovery and wrongly throttling the congestion window upon random packet loss can result in extremely poor TCP performance under typical WWAN operating conditions, as we have observed in our experiments with the CDPD network.

In recent years, there have been several proposals to optimize TCP for wireless networks. These proposals typically have had three flavors: (a) improving reliability at the link layer [6], (b) providing TCP-aware smarts in the base station[4], and (c) splitting the TCP connection into two parts [3, 13] - one in the backbone network and one over the wireless link running a specialized protocol [13]. Most of these solutions focus almost exclusively on making the loss characteristics of the wireless link transparent to the sender, while preserving the congestion control mechanisms of TCP. As we will discuss in Section 2, contemporary solutions address only a part of the problem of WWAN environments and leave a lot of scope for improvement.

In this paper, our goal is to identify the fundamental causes of performance degradation of TCP in commercial WWAN networks, and design a transport protocol that effectively addresses these causes. We seek to design a deployable, robust, fair, efficient, and reliable transport protocol for commercial WWAN networks. To this end, we present the *Wireless Transmission Control Protocol (WTCP)*, which enunciates the following three key principles:

1. WTCP uses purely *end-to-end mechanisms*, thereby eliminating the need for transport-level support from the network infrastructure, in contrast to related work [6, 4, 3, 13]. We choose end-to-end mechanisms for three reasons. First, we believe that end-to-end mech-

anisms are necessary to effectively address the problems of WWAN environments. Second, the connection end-points in most of the current WWAN deployment scenarios are a mobile host and a dedicated proxy - both of which are already wireless/mobility aware. It is thus feasible to propose end-to-end transport layer solutions for WWAN environments without having to change the TCP implementations on all stationary hosts. Third, WWAN network providers are typically unwilling to introduce TCP-aware smarts in the base station or mobile switching station for reasons discussed in Section 2.

2. WTCP uses *rate-based* rather than window-based transmission control. As a result, WTCP shapes its data traffic, never allows a burst of packet transmissions, and is also less prone to the unfairness observed in TCP when competing connections have different round-trip times [9]. Furthermore, it is the receiver that adaptively computes the desired transmission rate based only on the characteristics of the data path. Consequently, WTCP is relatively insensitive to problems in the ACK path and can handle asymmetric channels well.
3. WTCP uses the ratio of the *inter-packet separation* at the receiver and the inter-packet separation at the sender (recall that the sender shapes data traffic) as the primary metric for rate control rather than using packet loss and retransmit timeouts. As a result, WTCP reduces the effect of non-congestion related packet losses on the computation of transmission rate. Additionally, WTCP is less prone to large round-trip-time variations as the congestion window size increases, which is inevitable in standard TCP-variants due to the low-bandwidth nature of the channel.

We currently have a user-level implementation of WTCP that we have evaluated over the CDPD network, and a simulation of WTCP in the ns-2 simulator. Initial performance results indicate that WTCP significantly improves on the state-of-the-art solutions for wireless transport.

The rest of the paper is organized as follows. Section 2 identifies the typical characteristics of the WWAN environment, and how they adversely affect the performance of TCP. Section 3 describes the key design principles and mechanisms of WTCP. Section 4 presents the algorithms for congestion control and reliability for WTCP. Section 5 evaluates the performance of the WTCP implementation over CDPD and simulation in ns-2. Section 6 compares WTCP to related work. Section 7 concludes this paper.

2 Characteristics of WWAN Environments

At a high level, most WWAN environments exhibit similar characteristics: the bandwidth is low and varying, the latency is high and varying, blackouts exceeding 10 seconds occur occasionally, and for unreliable WWAN networks, non-congestion related packet loss may be significant when the user is traveling at moderate to high speeds.

We have used CDPD as the evaluation platform for WTCP. Two characteristics of CDPD are germane to our discussion: the effective throughput of a CDPD channel typically does not exceed 12Kbps, and the majority ($\geq 75\%$) of the end-to-end latency is incurred in the CDPD part of the network between the mobile switching station and the mobile host. We will revisit the impact of the latter point when we discuss

wireless transport protocols that rely on TCP-aware smarts at the base station [4].

WWAN wireless networks in general, and CDPD networks in particular, typically exhibit the following characteristics.

1. *Non-congestion related packet loss*: Even though CDPD uses Reed-Solomon(63,47) error-correcting code, we have measured non-congestion related errors ranging from 0% to 10% depending on the speed of mobility (measured over a range of 0 - 55mph), location of the user vis-a-vis the base station, and co-channel interference. TCP assumes that all packet losses result from congestion. A 5% non-congestion related packet loss can thus significantly degrade the performance of TCP.
2. *Very low bandwidth*: Between 1 to 30 users may share a single CDPD channel of raw capacity 19.2Kbps. For RAM, the channel is 8Kbps, while for Ardis, the channel may be 4.8Kbps or 19.2Kbps. Due to the extremely low bandwidth, the delay-bandwidth product of a connection is small (typically 2 or 3 packets). This can affect the congestion control and fast retransmit mechanisms of TCP adversely. TCP sometimes observes artificially larger congestion windows as a result of deep buffering in the CDPD network. While this allows a connection to pump in more packets into the network, sending back to back packets on a slow link can artificially increase the round trip time and adversely affects TCP performance in case of a timeout.
3. *Large round trip time and variance in round trip time*: In CDPD, we have observed typical round trip times between 800ms to 4 sec. A large fraction of this time is due to transmission on the wireless link (e.g. transmitting a 512 byte packet at 12Kbps takes 300ms), and over 75% of the latency is typically incurred in the mobile switching station \leftrightarrow mobile host segment of the connection.

TCP sets the retransmission timeout to be the sum of the average round-trip time and twice the mean deviation of the round-trip time ($RTO \leftarrow \overline{rtt} + 2 \times \Delta(rtt)$). This may result in very large RTOs (e.g. 32 sec) in CDPD because of two reasons: (a) the rtt and $\Delta(rtt)$ are inherently large, and (b) ACKs from the receiver get bunched (see below), and since ACKs clock data packets in TCP, data packets are sent out in bursts, which further increases the mean and deviation of rtt . Thus, timeouts affect TCP performance severely in CDPD.

4. *Asymmetric Channel - Bunching of ACKs*: CDPD uses DSMA/CD [2] for contention resolution in the channel. Contentions among mobile users for the up-link channel are resolved by binary exponential backoff. Consequently, CDPD suffers from the well known 'capture syndrome' of binary exponential backoff [11], in which a highly loaded shared medium ends up bursting the queued packet transmissions of each contending host in turn. RAM also suffers from the same problem. For the common case of downlink data transmission, ACKs from the mobile to the backbone host get bunched. This further skews the round trip time computation, and also causes the sender to burst out packets as mentioned above.

5. *Occasional blackouts*: Prolonged fades, sudden degradation in signal quality such as traveling through a tunnel or between non-overlapping base stations, and temporary lack of available channels (when cellular phone calls are occupying the channels) can cause blackouts lasting 10 seconds or more, and result in the back-to-back loss of a sequence of packets. Traveling at 55mph, we observed several blackouts ranging from 10 seconds to 10 minutes during the course of a day.

2.1 High Level Architectural Trade-offs in WWAN Environments

Before going into the details of the WTCP design, we need to step back and discuss the trade-offs between using end-to-end mechanisms for wireless TCP versus using smart mechanisms in the network in order to assist wireless-unaware TCP end-points. Related work on wireless optimizations for TCP has typically argued against end-to-end mechanisms on the grounds that it is impractical to change the protocol stacks of all stationary hosts merely to accommodate mobile hosts. Thus, most of the previous work has focused on making the lossy nature of the wireless link transparent to the stationary end host by introducing smarts at the base station via one of three mechanisms: reliable link layers [6], TCP-aware ‘snoop’ mechanisms [4], or splitting the connection into two (wireline and wireless) distinct components [3, 13].

Link layer retransmission [6] works well when the latency over the wireless link is small compared to the coarse grain TCP timer. In the ideal case, the link layer retransmissions are not expected to significantly interfere with the end-to-end *rtt* computations or congestion control mechanisms - except to eliminate random channel loss. In WWAN networks, it is the transmission time over the wireless network that constitutes the bulk of the observed end-to-end latency. Consequently, providing only a reliable link layer abstraction and keeping TCP unchanged at the end hosts simply will not work because link-layer retransmissions will interfere with the reliability and congestion control mechanisms of TCP.

The Snoop protocol [4] instantiates TCP-aware smarts at the base station (or mobile switching station) in order to eliminate the problem of false fast retransmits or slow starts due to random packet loss over the wireless channel. As described in Section 2, this approach also assumes that the transmission time over the wireless link is significantly smaller than the coarse grain TCP timer and round trip time. In wireless LAN environments, snoop works well. However in WWAN environments, snoop does not work well because of two reasons: (a) it exacerbates the problem of large and varying round trip times by suppressing duplicate ACKs, and (b) duplicate retransmissions may be initiated by both the snoop agent and the end host (which may observe a timeout) because of comparable timeout values at the two entities. In fact, we have observed that snooping may possibly degrade the performance of TCP when the latency over the wireless link dominates the round trip time. Additionally, snoop does not work well with some flavors of TCP such as TCP-Vegas. In summary, we believe that snoop works well in the environment for which it was designed, but it does not work well in the WWAN environment.

Indirect TCP protocols [3, 13] break the TCP connection at the base station, and maintain two separate connections - one over the wireline network and one over the wireless network. I-TCP violates the fundamental end-to-end guarantees of TCP by splitting the connection. Note that the

connection split must happen at the base station (or mobile switching station) serving the mobile host, and the connection state must be moved across base stations upon hand-off. The wireless component of I-TCP is quite simplistic and does not address several of the issues listed above. As a practical matter, the I-TCP architecture may not be feasible for WWANs because it requires significant infrastructure support and maintenance of connection state from the WWAN network, which is autonomously managed and may not even understand TCP/IP internally (e.g. RAM). It is important to note that I-TCP calls for splitting each transport connection transparently to the TCP end-points. It is thus quite distinct from the WWAN deployment scenario of a mobile host connecting to a dedicated stationary proxy in the backbone.

To summarize, we believe that previous approaches that seek to hide the problems of WWAN networks from the TCP at the end host by adding TCP-aware smarts in the mobile switching station are not applicable to WWAN networks for two reasons: (a) such approaches require the base station to maintain significant state, understand TCP/IP, and are often tuned to specific flavors of TCP, and (b) the fact that the latency between the mobile switch station and the mobile host is the dominant component of a large and varying round trip time makes such approaches less effective. The bottomline is that for WWAN environments, both end-points must cooperatively address the issues unique to the environment. Moreover, it is desirable to eliminate network-level smarts because the base stations are owned by an autonomous entity that may not even be running IP internally. We believe that the key issues that need to be addressed - the non-congestion related packet loss, large and highly varying latency, asymmetry in data/ACK channel behavior - can be effectively solved with the end-to-end mechanisms proposed in this paper. Of course, the penalty for using the end-to-end mechanism is that the remote end-point in the backbone must also change. Luckily, the nature of the WWAN environment and the current deployment pattern already supports the common case of WWAN users typically connecting through a dedicated proxy server on the backbone.

3 The WTCP Approach

Any reliable transport protocol must provide the following functions: (a) connection management, (b) congestion control, (c) flow control, and (d) reliability. The connection management and flow control in WTCP are similar to the standard TCP mechanisms and are not described in this paper. We now focus on the key design choices in WTCP for congestion control and reliability.

3.1 Congestion Control

The key aspects of congestion control in WTCP are that it is rate based, uses inter-packet delay as the primary mechanism to determine rate adaptation, performs the rate adaptation computations at the receiver, predicts the cause of packet loss and reacts accordingly, and varies the granularity of rate increase/decrease depending on the type of congestion observed. Additionally, WTCP also tailors its startup behavior to work well for short-lived flows. We describe the design aspects of WTCP congestion control in more detail below.

1. *Rate-based transmission control*: As we mentioned in Section 2 for the common case of bulk data trans-

fer from the backbone host to the mobile host, ACKs are often bunched together on the return path to the sender because of the nature of channel arbitration in CDPD. With the window-based self-clocking mechanism of TCP, this results in the sender bursting back-to-back data packets, which skews round trip time computations, causes more bursty queuing at the base station, and consequently more packet drops. WTCP alleviates these problems by using a rate based scheme that does not use ACKs for self-clocking. Adopting a rate-based approach does involve explicit clocking and shaping the traffic according to the current transmission rate of the connection; however, the rates are typically small enough that coarse grain timers ($O(100ms)$) are sufficient to perform the clocking effectively.

2. *Inter-packet delay as the main mechanism for transmission control:* We have observed that monitoring the average inter-packet delay at the receiver provides a fairly accurate measure of the available channel rate. Specifically, the ratio of the sender's inter-packet delay and the receiver's observed average inter-packet delay provides a responsive metric to determine whether the transmission rate is in steady state, can be increased, or should be decreased. Thus, when the network is uncongested or has incipient congestion, reacting to changes in the inter-packet delay ratio serves to keep the network uncongested, and significantly reduces the number of congestion-related packet losses. WTCP thus uses this mechanism as the primary transmission rate control mechanism, and essentially uses incipient congestion detection without waiting to lose packets before throttling down the sending rate.
3. *Distinguishing the cause of packet loss and adjusting transmission rate accordingly:* While inter-packet delay is the main mechanism for dealing with incipient congestion, if the network suddenly moves from uncongested to congested state (e.g. due to a sudden influx of new connections or sudden decrease in available resources), then packets are dropped due to congestion. In WTCP, the receiver maintains a history of packet losses when the network is uncongested, and computes the expected average and deviation in the number of non-congestion-related packet losses over a time window. Based on this information, it predicts if an observed packet loss is congestion related or not. If the loss is predicted to be due to congestion, then the sending rate is throttled down. Otherwise, the sending rate is adjusted according to mechanism 2 above (to deal with the uncongested or incipient congestion case). We typically observe that congestion-related throttling does not occur often since most of the random losses conform to the pattern generated from the history. However, during blackout periods, this mechanism erroneously assigns blackout-related packet losses to congestion. This causes the receiver to throttle down the sending rate (though the ACKs may not get through to the sender if the blackout is two-way). Once the receiver determines that a blackout has occurred (the sender goes into blackout mode as described in point 6 below), it resets the sending rate to the old (pre-congestion) value upon recovery from blackout.
4. *Performing transmission control computations at the receiver:* In WTCP, the receiver performs the rate control mechanisms described above, and computes the

new transmission rate of the sender. With each data packet, the sender transmits its current inter-packet separation. Based on local state and the state in the packet, the receiver has all the information it needs to update the transmission rate, which it then sends back in its acknowledgment. Having the receiver perform the rate computations eliminates the effect of delay variations and losses in the ACK path. Even if ACKs get bunched, delayed, or lost, the transmission rate is not altered. WTCP can thus deal with asymmetric channels effectively.

5. *Variable granularity rate adjustment:* TCP uses the well known linear-increase-multiplicative-decrease policy for adjusting its congestion window. While LIMD is stable and asymptotically converges to fair channel allocation, the efficiency of the LIMD algorithm is a function of how severely the decrease is performed. TCP reduces its congestion window by half upon observing a packet loss. In WTCP, we seek to detect incipient congestion and react to it early on in the common case. The goal of WTCP is to decrease the transmission rate multiplicatively in order to ensure fairness, less aggressively when reacting to incipient congestion in order to improve efficiency, and more aggressively when reacting to real congestion in order to reduce packet loss and alleviate congestion quickly. In order to achieve these goals, WTCP maintains a history of transmission increase/decrease in the recent past. If the receiver is required to perform transmission decrease multiple times in quick succession, it starts to decrease its transmission rate more aggressively. If the receiver observes a congestion based packet loss, it halves its rate. As a result of this approach, incipient congestion is handled by a gentle decrease of the transmission rate, but severe congestion is handled by an aggressive decrease in transmission rate.
6. *Startup behavior:* Since round trips are large in WWAN environments, and since some data transmissions may be short-lived, WTCP attempts to compute the appropriate transmission rate for a connection immediately upon startup rather than going through slow start. WTCP uses the 'packet-pair' approach [8], wherein it sends two back-to-back packets of maximum segment size (MSS) and computes their inter-packet delay during connection establishment. This serves as an approximate estimate for the sending rate.

This mechanism can also be used to recover from blackouts, which occur when the connection may experience back-to-back losses for extended periods of time due to poor channel conditions or lack of available channels. We have also considered an alternative mechanism for blackout recovery, in which the sender starts to 'probe' the receiver (instead of sending packet-pairs) when it does not receive an ACK for a threshold period of time. The difference in the two cases is whether a successful packet handshake after the blackout reverts to the old transmission rate or computes a new transmission rate according to the new channel conditions. A comparison of these schemes is ongoing work.

3.2 Reliability

The key aspects of reliability in WTCP are that it uses selective acknowledgments, it does not use retransmit timeouts,

and that it tunes the frequency of sending acknowledgments to the dynamic network conditions. We describe these aspects below:

1. *Selective Acknowledgments*: As noted in related work, selective acknowledgments are very useful in TCP [10]. WTCP uses selective acknowledgments for ensuring reliability. The receiver periodically sends ACKs at a frequency tuned by the sender (see below), containing the cumulative and selective ACK. By inspecting the ACK, the sender can detect a hole in the receiver's sequence of received packets. By comparing the state contained in the ACK with local state stored with the last (re)transmission for each unacked packet, the sender can determine if this last (re)transmission was lost, or could still be in transit. Thus selective acknowledgment allows the sender to retransmit only lost packets.
2. *No retransmit timeouts*: As we have observed in Section 2, it is exceedingly difficult to maintain a reliable estimate of the retransmit timeout. In fact, many of the performance problems observed in various TCP flavors are caused by erroneous RTO estimation. WTCP does not use RTOs. Instead, it *modifies the SACK algorithm* to achieve reliable transmission without RTOs. This mechanism is described in Section 4, and is an important aspect of WTCP.
3. *Controlling ACK frequency*: ACKs carry both reliability and transmission control information, and the sender must receive ACKs periodically in order to react to the new transmission rate, and perform flow control. The sender tunes the desired ACK frequency (and notifies the receiver in the data packet) such that it expects to receive at least one ACK in a threshold period of time (e.g. 5 seconds). If the sender does not receive an ACK in the threshold period, it goes into blackout mode. The tuning of the ACK frequency is governed by several factors: (a) observed ACK loss at the sender, (b) half-duplex or full-duplex nature of the WWAN channel, and (c) average and deviation in the inter-ACK separation observed at the sender. Note that a receiver may also voluntarily generate a SACK immediately upon observing a hole in the packet sequence.

4 The WTCP Algorithm

In WTCP, *rate control and reliability are decoupled*. The receiver computes the desired sending rate via its rate control mechanisms, and notifies this rate to the sender in the ACK packets. ACKs thus carry both reliability information (SACK) and rate control information. The sender monitors the reception of ACKs, and adjusts its rate accordingly. It also monitors the ACKs to tune the ACKing frequency, which it then notifies to the receiver in future data packets. If the sender does not receive an ACK for a threshold period of time, it goes into blackout mode and periodically sends probe packets to elicit ACKs from the receiver and recover from the blackout. The probe packet mechanism is also used for loss recovery, eliminating the need for timeout-based retransmission in WTCP.

We now describe the parameters contained in the packets of WTCP, and then describe the key functions of WTCP. There are three types of packets: Data, ACK, and Probe. The contents of these packets are as follows:

Data: $\langle rel_seq_num, cctrl_seq_num, hack_seq_num, current_rate, ack_frequency, packet_size, data \rangle$
ACK: $\langle ack_seq_num, hcctrl_seq_num, updated_rate, CACK, SACK \rangle$
Probe: $\langle cctrl_seq_num, optional_data \rangle$

The 'reliability sequence number' (rel_seq_num) in the data packet is the byte sequence number as in TCP, and is used for sequencing and reliable delivery of packets. The 'congestion control sequence number' ($cctrl_seq_num$) is a monotonically increasing packet sequence number that distinguishes multiple transmissions of the same packet. The ACK packet contains the CACK and SACK that acknowledge maximal sequences of contiguous data packets that have been received at the receiver, and also the highest congestion control sequence number ($hcctrl_seq_num$) seen thus far by the receiver. This is used by the sender in a variant of the SACK algorithm to eliminate retransmit timeouts.

ACKs contain the desired sending rate and are sequenced by the ack sequence number (ack_seq_num), which is monotonically increasing for every ack sent from the receiver. Each data packet from the sender contains the highest ack sequence number ($hack_seq_num$) that it has seen so far, and this information is used at the receiver in rate control to determine if a sender has already reacted to its previous rate adjustment. Other parameters in the data packet include the current sending rate (inverse of sending inter-packet delay), the desired ACK frequency, packet size, and data. The probe packet contains a congestion control sequence number ($cctrl_seq_num$) and optional data. As in the case of data packets, sending a probe packet also increments the $cctrl_seq_num$ by 1.

We now describe the rate control and reliability algorithms in WTCP.

4.1 Rate Control

In order to perform rate control, the receiver maintains the following state: (a) long-term (avg_ratio) and short-term running averages ($svag_ratio$) of the ratio of the actual rate of the sender to the observed rate at the receiver, (b) the maximum expected number of random losses in a given window of packets, and (c) a decrease variable (δ) that indicates the percentage decrease in the desired sending rate if the sending rate needs to be throttled.

At any time, the receiver can be in one of three states:

- *increase*: the long-term and short-term ratios of the observed sending rate at the receiver to the actual sending rate at the sender are both above certain thresholds (i.e. ($avg_ratio > \alpha_+$) and ($svag_ratio > \beta_+$)).
- *decrease*: the long-term and short-term ratios of the observed sending rate at the receiver to the actual sending rate at the sender are both below certain thresholds (i.e. ($avg_ratio < \alpha_-$) or ($svag_ratio < \beta_-$)).
- *maintain*: otherwise.

When the receiver receives a data packet, it checks to see if the packet has the next expected congestion control sequence number. If the packet is received in sequence, avg_ratio and $svag_ratio$ are updated. Then the receiver transitions to the *increase*, *maintain*, or *decrease* state, depending on the ratios. For increase, the short term average captures the effect of increase in available bandwidth, but to avoid a false increase, the long term average is also required to be above a certain threshold (α_+), which is very close to

```

decrease_sending_rate()
1  decrease sending rate by  $\delta$  %
2   $\delta \leftarrow \text{MAX}(2 \times \delta, 50)$ 
3  current_phase  $\leftarrow$  decrease

process_packet()
1  if (loss_observed was reset more than  $T_{loss}$  time units ago)
2    reset loss_observed
3  if (packet contains data and
4    previous data packet was missing)
5    update loss_observed
6    if (current_phase = maintain)
7      update loss_expected
8    if (loss_observed > loss_expected)
9      decrease sending rate by 50%
10      $\delta \leftarrow 20$ 
11     current_phase  $\leftarrow$  decrease
12 else if (packet contains data) /* not a probe or an ack */
13   update lavg_ratio
14   update savg_ratio
15   if (sender has started using the latest computed rate)
16     if (lavg_ratio  $\geq$   $\alpha_+$  and savg_ratio  $\geq$   $\beta_+$ )
17       increase sending rate by a constant
18        $\delta \leftarrow 20$ 
19       current_phase  $\leftarrow$  increase
20     else if (lavg_ratio <  $\alpha_-$  or savg_ratio <  $\beta_-$ )
21       decrease_sending_rate()
22     else /* Maintain */
23       set the new rate to be same as the old rate
24        $\delta \leftarrow 20$ 
25       current_phase  $\leftarrow$  maintain
26   else if (last action was a speedup)
27     if (lavg_ratio <  $\alpha_-$  and savg_ratio <  $\beta_-$ )
28       decrease_sending_rate()
29   send ack

```

Figure 1: Rate Control for WTCP at the receiver

one. So only when both these conditions are satisfied, can we safely increase the sending rate. Similarly for decrease, the short term average captures the effect of reduction in available bandwidth. But to avoid the case where the short term ratios remain between β_+ and β_- , though the network is in a congested state, the long term average used is required to be lower than a certain threshold (α_-) which is very close to one. The combination of both these averages makes the increase as well as the decrease steps judicious.

The main components of the congestion control algorithm are described below, with references to the pseudocode in Figure 1.

1. *RTT computation:* As discussed earlier, retransmission timeouts are not used in WTCP, however a rough estimate of RTT is maintained because, when the WTCP sender goes into the blackout mode, it sends probe packets at an interval of $2 \times RTT_{est}$. To maintain the RTT estimate, when the sender (re)transmits a packet, it records the transmission time of the packet. The first time that the packet is acknowledged, the sender records the round trip time for the packet. Our RTT estimates are reliable even for retransmitted packets, because the reliability mechanisms of WTCP, retransmits a packet only after determining that the last transmission was lost.
2. *Using loss profile:* WTCP keeps track of the random loss percentage which is computed/updated during the *maintain* phase of the algorithm, as in that phase, losses are least likely to be because of congestion. When a non-consecutive (in terms of congestion control sequence number) data packet is received, the observed loss percentage is updated (Line 4) based on the gap in the congestion control sequence numbers at the receiver. If the receiver is in currently the *maintain* phase, then the expected loss percentage is also updated (Lines 5-6).
If the observed loss exceeds the worst case expected loss, then the receiver predicts that there has been some congestion loss, halves the desired sending rate, and enters the decrease state.
3. *Increase Phase:* If the received packet is the next expected packet according to the congestion control sequence number, the long term and the short term averages of the ratio of receiving to sending rate are compared with the constants α_+ and β_+ respectively. If the receiving rate is greater than the sending rate or slightly less than the sender's rate, WTCP tries to do probing (Lines 15-18). If the increase phase is entered, the decrease variable δ is reset to 20.

4. *Decrease Phase:* WTCP decreases the rate by half on detection of a congestion loss. In other cases, when it decides to decrease based on the long term and the short term averages, the decrease variable δ is used. The initial value of δ is 20. As typically most of the decreases are because of probing, a reduction by 20%, improves the channel utilization to 90% (compared to 75% for TCP Reno). To take care of an actual congestion loss, WTCP doubles the decrease variable for consecutive decreases. Thus, it takes at most one extra RTT, compared to TCP Reno, to drop to 50% (decrease by 20%, followed by a decrease by 40%, as δ is doubled after a decrease, is approximately equivalent to a single drop of 50%) for a real congestion.

The algorithm enters the decrease phase from three different places (Lines 7, 20 and 27). As described earlier, the loss profile can indicate that a loss is because of congestion, in which case the decrease phase is entered (Lines 8-10). The decrease phase is also entered (Lines 19-20) if the sender has started using the latest updated rate and if the long term and the short term averages are below the thresholds α_- and β_- respectively. But if the receiver detects congestion after an increase was sent to the sender, then the receiver doesn't wait for the sender to start using the new rate, instead it sends a slow down right away (Lines 25-27). The idea behind comparison of the average ratios with α_- and β_- (Lines 19 and 26) is to enter the decrease phase only when WTCP notices that the receiving rate has fallen below the sending rate by a significant amount. Since the short term average (receiving rate to sending rate) might be affected due to transient network behavior, both the long term and the short term averages of the ratios are taken into account.

5. *Maintain Phase*: This phase is entered when the long term and the short term average ratios indicate that no increase or decrease needs to be done (Lines 21-24). The decrease variable δ is reset to 20 and the current_phase is set to *maintain*.

4.2 Reliability

As described earlier (Section 2), RTT estimates are inaccurate primarily due to the low bandwidth nature of the wireless links. Hence, WTCP doesn't have a notion of a retransmission timer. WTCP depends on the following two techniques for retransmission of packets.

1. *SACK Algorithm*: The important steps in processing an acknowledgment are shown in Figure 2 as a pseudocode. For every transmitted packet the sender keeps track of the congestion control sequence number of the packet which was used to send it. Lines 2 through 8 are executed for all unacked packets. If the packet is being CACKed or SACKed, then it is removed from the unacked list (Lines 2-5). *p.rel_seq_num* represents the reliability sequence number of the packet. A packet is marked for retransmission when the sender recognizes that the receiver has seen a packet with a larger congestion control sequence number. Packets marked ready for retransmission have a higher transmission priority compared to the new data packets.
2. *Probe Packets*: While the SACK algorithm described above effectively detects 'holes', it will not recover the loss of a packet unless another packet that was subsequently transmitted has been successfully received at the receiver. Consequently, it is possible for the sender to lose a sequence of packets, and then run out of more packets to send. Since there is no concept of a retransmission timer, lost packets would never be recovered.

In order to solve this problem, at the end of a run of packets, or upon not receiving an acknowledgment for a threshold period of time, the sender sends out periodic probe packets with increasing congestion control sequence numbers. The receiver responds to a probe packet with an acknowledgment, and now the

hcctrl_seq_num field in the ACK contains a congestion control sequence number that is larger than that of any transmitted packet. As a result, the SACK algorithm described above kicks in for the detection and retransmission of lost packets (Figure 2, Lines 6-8).

5 Performance Evaluation

5.1 Experiments on the CDPD network

In this section we first provide an overview of the characteristics that we observed with a CDPD network and then compare our user level implementation of WTCP against an off-the-shelf implementation of TCP-Reno available with Windows 95C and Linux 2.0.33. For all our experiments, we used a Sierra Wireless AirCard for CDPD access. The sender was a 233 MHz Pentium PC running Linux 2.0.33 (the fixed host) and the receiver was a 166 MHz Pentium Winbook laptop running Windows 95C (the mobile host). The CDPD experiments were performed in Chicago, IL and Bloomington, IL at speeds of 0-55 mph at different times of the day. Due to space constraints, we only present a small set of results obtained while traveling at 55 mph. Also, for all results reported here, we have used a fixed packet size of 512 bytes.

5.1.1 CDPD network characteristics

While transmitting UDP packets at 1 sec. intervals over 30 minutes, the observed short term error rate varied between 1% to 15% with a long term average of around 4%. To study the effect of packet bursts on round trip time, we sent packets from the mobile host to the echo port of the fixed host at different burst sizes. The observed RTT increased linearly within a burst. The typical value of RTT for a packet was about 1.8 sec. Figure 3(a) shows results from a experiment with a burst size of 8. This shows that bursting data packets would result in high RTT values and high variance in RTTs. Figure 3(b) shows the ratio of sender to receiver rate, its short term average (50% weightage on a new sample), long term average (12.5% weightage on a new sample), and the mean deviation of the long term average (12.5% weightage on a new sample). These were obtained by sending packets from the fixed host with a send separation of 1sec. The graph shows that the short term average is able to capture current network behavior, and the long term average remains close to 1.0 in steady state.

5.1.2 TCP vs. WTCP

To compare the performances of TCP Reno and WTCP, we sent 100 KB of data from the fixed host. On an average with a packet loss rate of about 4%, TCP takes 134 seconds to send 100KB of data from the fixed host to the mobile. Note that without slow-start WTCP takes less than 100 seconds at an average for sending the same amount of data. The total time taken by WTCP does not vary with varying loss percentage, as the sending rate in WTCP is not affected by random losses. In Figures 4 and 5 we present two TCP and WTCP performance results, which were obtained while moving at a speed of 55 mph.

In Figure 4(a), there are 8 random losses. TCP stalls at periods as it is waiting for reception of acks to move/increase its congestion window. There are two retransmission timeouts, one at 35 sec (RTO 30 sec) and another at 128 sec (RTO 17 sec). For all other losses, TCP is able to recover using fast retransmit. The large retransmission timeout at

```

process_ack(acknowledgment a)
1  foreach packet p in the unacked list
2      if ( $p.rel\_seq\_num < a.CACK$ ) /* The cumulative ack is acking the packet */
3          remove packet from unacked list
4      else if ( $p$  is SACKed) /* The selective ack is acking the packet */
5          remove packet from unacked list
6      else if ( $p.ready\_for\_retx = FALSE$ ) /* Packet hasn't been marked */
7          and ( $p.ctrl\_seq\_num < a.hcctl\_seq\_num$ ) /* Ack corresponding to a later packet */
8           $p.ready\_for\_retx \leftarrow TRUE$ 
9   $hack\_num\_rcvd \leftarrow \max(a.ack\_seq\_num, hack\_num\_rcvd)$  /* Update the highest ack number received */

```

Figure 2: *The Reliability Algorithm for WTCP at the sender*

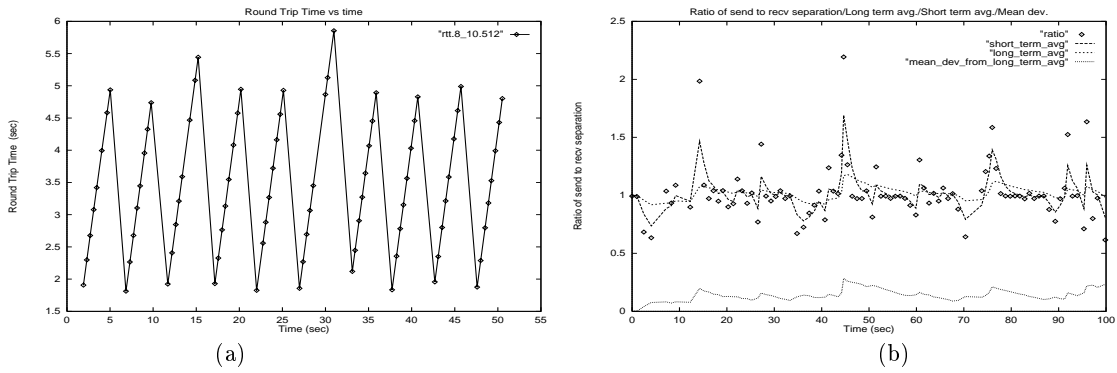


Figure 3: (a) RTT for burst size of 8 with packet size of 512 bytes. (b) Ratio of send to rcv separation for packet sizes of 512 bytes sent at 1 sec. interval. Also shown are the short term average, long term average and the mean deviation for the long term average.

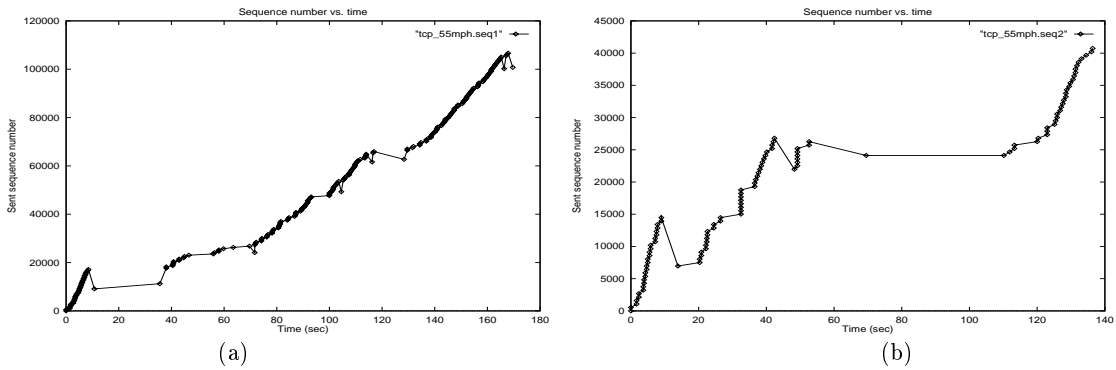


Figure 4: Performance of TCP on CDPD while moving at 55 mph. A blackout period is captured in Figure (b)

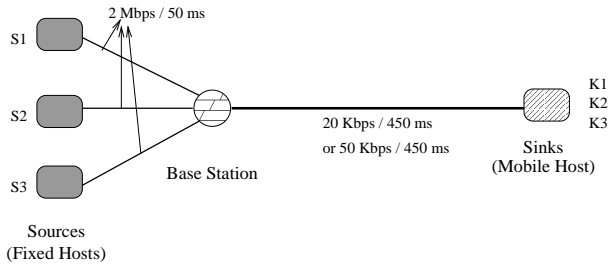


Figure 6: Network topology

approximately 35 sec. is because of the slow start mechanism which aggressively fills up the buffers, resulting in large RTTs and correspondingly large RTOs. In Figure 4 (b), there are two sequences of losses one from 13 sec to 26 sec and the other from 42 sec to 120 sec. During these blackout periods, TCP reacts using the RTO mechanism and backs off its RTO value resulting in poor performance.

In Figure 5 (a), there are 11 losses. At around 70 sec, WTCP enters into blackout and stops transmission for a few seconds, till it starts getting acknowledgments from the receiver. But overall, the rate of transmission is not effected by the random losses. In Figure 5 (b), there are 3 random losses and a blackout period of about 60 seconds. As soon as the blackout period is over, WTCP starts sending at the rate it was using before going into blackout.

A comparison of Figures 4(a) and 5(a) shows that for approximately the same number of losses, TCP takes 170 seconds for sending 100KB of data whereas WTCP takes 130 seconds (inclusive of a short blackout). Similarly for Figures 4(b) and 5(b), we find that TCP takes 140 seconds to send 40 KB of data, whereas in the same amount of time, WTCP is able to send about 80 KB of data. As TCP does not distinguish between congestion and random losses and has no mechanism to handle blackouts, these results were expected of TCP. Also, as WTCP has smarts to distinguish between congestion and random losses, and also recognizes blackouts, expectedly, it performs much better than TCP on the real CDPD network.

5.2 Simulations

In this section we compare the performance of the congestion control algorithms in WTCP, TCP-Newreno[7], TCP-Vegas[5] and TCP with snoop[4] via simulations using the ns-2 simulator². The network topology used in the simulations is shown in Figure 6. Though this topology does not mimic all the characteristics of any specific wireless wide area network, it gives a clear picture of the performance of various transport protocols over channels with low bandwidth and significant error rate. While networks such as CDPD currently offer only 19.2Kbps raw bandwidth, it is expected to grow by two or three times in the near future [1]. So results have been presented using 20Kbps and 50Kbps bandwidth for the wireless channel. For the simulations we have used an exponential error model with mean error rates from 3 to 8%. All the simulations presented in this section use a packet size of 500 bytes.

In the rest of this section, we present three different sets of results. For the first set of results we use a single flow to evaluate the performance of TCP-Newreno, TCP-Vegas and

WTCP with different bandwidths and error rates. Then we illustrate the fairness of WTCP using two flows. Finally, we compare WTCP with different versions of TCP using the Snoop protocol. As congestion control is primarily concerned with the total number of packets successfully sent and not the specific sequence number of the packet, we have used the total number of packets successfully transmitted as the primary metric to compare their performances. We have presented here only a small subset of the tests that we have performed. More detailed simulation results can be obtained from [12].

5.2.1 Single flow

The results for this scenario shown in Figures 7 (a), (b) and (c), use a wireless bandwidth of 50Kbps, a queue size of 5 packets and error rates of 0%, 4% and 6% respectively. Figure 7 (a) demonstrates that WTCP provides the same throughput as TCP-Vegas when the channel is error free. The next two graphs show that the performance of Newreno and Vegas degrade significantly with increasing error rates. The figures indicate that for a 50Kbps channel with an error rate of 4%, WTCP provides about 100% improvement in performance over the other versions of TCP we have considered. Figures 8 (a), (b) and (c) use a wireless bandwidth of 20Kbps and error rates of 4%, 6% and 8% respectively. The difference in performance between WTCP and other mechanisms are less pronounced with a smaller bandwidth because for a single flow, the packets already in the queue keep the channel occupied for a larger period of time, thereby reducing the impact of the drastic reduction of the sender's congestion window. Also with high error rates, the maximum channel utilization that can be achieved by TCP is bounded by the error rates itself. The results in this section indicate that in wide area wireless networks WTCP clearly outperforms its wire-line end to end counterparts.

5.2.2 Multiple flows

In this subsection, we study the behavior of WTCP with multiple flows. For all the cases considered we have used a queue size of 50 packets and bandwidth of 50Kbps and an error rate of 4% for the wireless channel. In the scenario shown in Figure 9 (a), a new flow is introduced into the system when another existing flow is in steady state. The introduction of a new flow results in an increase in the receive packet separation for the old flow. As a result, the old flow increases its send packet separation thereby providing more room for the new flow. The figure indicates that the new flow manages to obtain the fair share of the wireless bandwidth in a short period of time. The parallel lines clearly show that the wireless bandwidth is equally divided between the new and the old flows. The result in this section indicates that the WTCP congestion control algorithm is fair.

5.2.3 Snoop Protocol

Figures 9 (a) and (b) correspond to the scenario in which a 50Kbps wireless channel was used with the snoop protocol being run at the base station. The queue size was set to 5 and the error rates were 4% and 6% respectively. Also the channel between the source and the base station was 1Mbps with a latency of 50ms and the latency between the base station and the sink being 200ms. Since the latency between the Snoop agent and the mobile host dominates the round trip time, the use of Snoop does not translate into

²The WTCP agents and the simulation scripts can be obtained from <http://timely.crhc.uiuc.edu/Projects/wtcp.html>

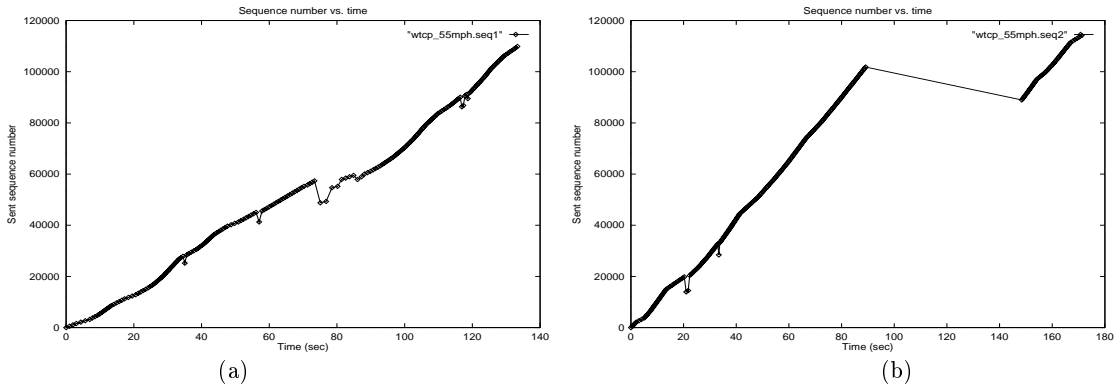


Figure 5: Performance of WTCP while moving at 55 mph. A blackout period was captured in Figure (b)

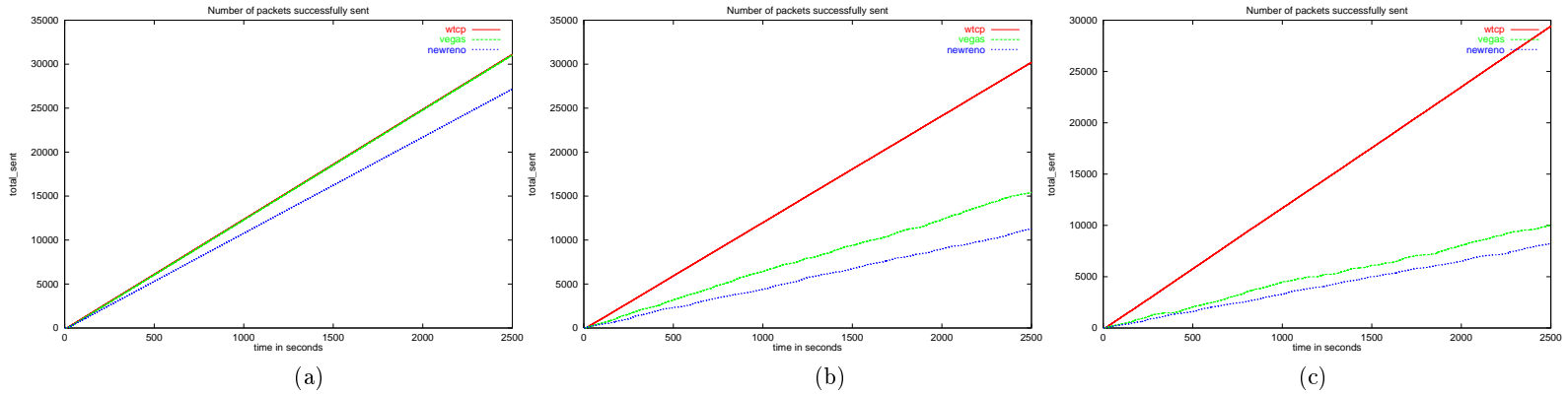


Figure 7: Total number of packets sent successfully with 50Kbps bandwidth and error rates (a) 0% (b) 4% and (c) 6%

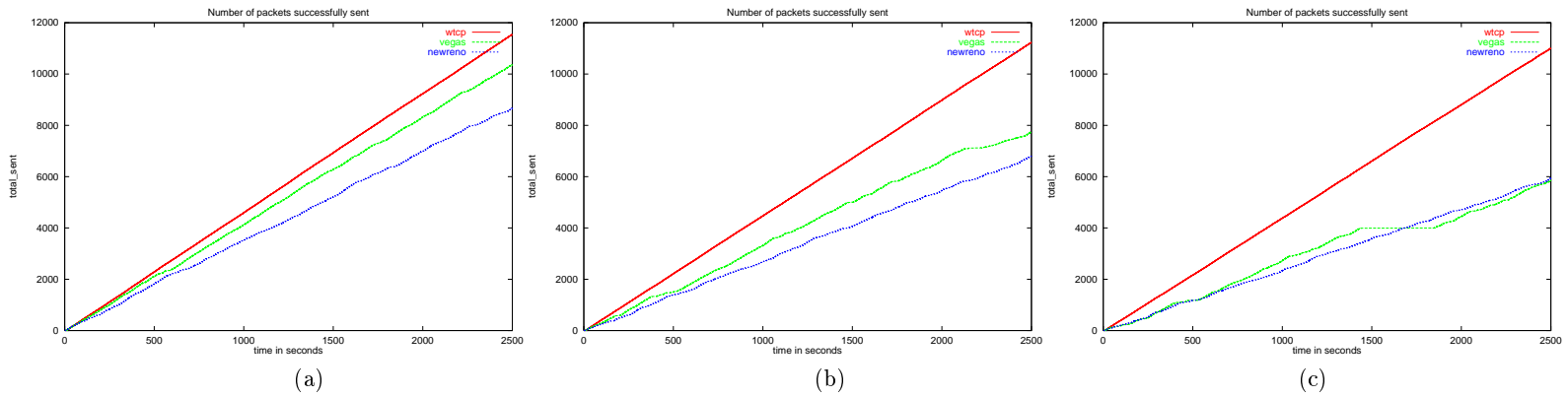


Figure 8: Total number of packets sent successfully with 20Kbps bandwidth and error rates (a) 4% (b) 6% and (c) 8%

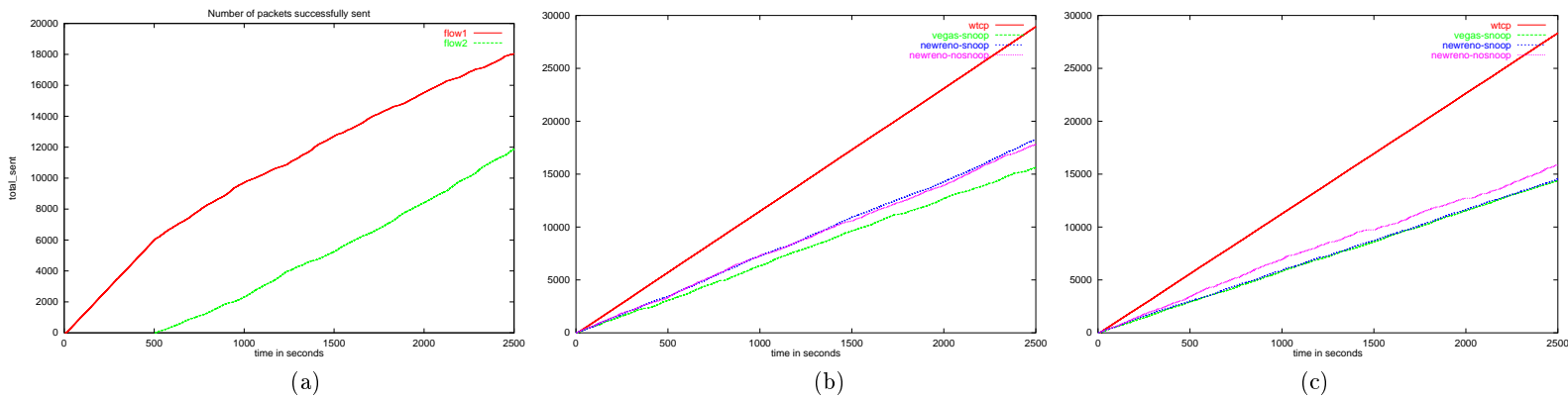


Figure 9: (a) WTCP with two flows starting at different times. (b) and (c) WTCP and TCP-snoop with 50Kbps bandwidth and error rates of 4% in (b) and 6% in (c)

an improvement in performance. Also the performance of TCP-Vegas deteriorates because Snoop interferes with the early congestion adaptation mechanisms. Figure 9 (b) considers the case where the error rate is 6%. In this case too, WTCP performs better than TCP-Newreno using Snoop. Note that the use of Snoop cannot improve the performance of TCP-Newreno beyond the case where there are no losses caused by channel errors. As we have seen in all the subsections above, even with losses, WTCP performs better than TCP-Newreno. Thus, Snoop is not very effective in hiding the losses in the wireless channel from the source, if the delay experienced after the point where the Snoop protocol executes is a large fraction of the total round trip delay. In other cases, as seen from the Figures 9 (a) and (b), it does not provide any improvement.

6 Conclusion

WTCP is rate-based and the rate adjustment is performed at the receiver; consequently, WTCP does not burst packets, overcomes the problems of inaccurate round-trip time computations, and handles asymmetric channels. WTCP uses inter-packet separation at the receiver as the primary metric for rate control with congestion-related loss detection as the backup mechanism; responding early to incipient congestion helps to keep the algorithm stable and at the same time, handling congestion-related losses causes WTCP to react correctly and fairly in worst-case scenarios of sudden congestion peaks. WTCP uses SACK and no retransmission timers for loss recovery; this enables efficient loss recovery without going into prolonged timeouts in the worst case. WTCP also provides for recovery from blackouts and good startup behavior for short-lived flows. In summary, WTCP handles most of the problems of WWAN, including those traditionally ignored by related work such as large and varying round trip times, and the significant fraction of the delay being incurred in the wireless segment.

We have not addressed several issues in the paper. First, we have not isolated the benefits of each of the mechanisms presented in this paper. Second, we have not compared alternate schemes for tuning ACK frequency and recovering from blackouts. Third, we have not evaluated the performance of WTCP for short-lived flows and more complex scenarios in detail. Fourth, we have not studied the effect of having deep buffers on performance in detail. Finally, we have not explored all the implications of our mechanisms on

fairness. These issues, along with a more robust publicly available user-level implementation of WTCP, are the foci of our ongoing work.

References

- [1] Wireless Data Forum. <http://www.wirelessdata.org/>.
- [2] John Agosta and Travis Russle. *CDPD: Cellular Digital Packet Data Standards and Technology*. McGraw Hill, New York, NY, 1997.
- [3] A. Bakre and B. R. Badrinath. I-TCP: Indirect TCP for mobile hosts. In *Proc. 15th International Conf. on Distributed Computing Systems (ICDCS)*, May 1995.
- [4] H. Balakrishnan, S. Seshan, E. Amir, and R. Katz. Improving TCP/IP performance over wireless networks. In *ACM MOBICOM*, November 1995.
- [5] L. Brakmo and L. Peterson. End to end congestion avoidance on a global internet. *IEEE JSAC*, 13(8), October 1995.
- [6] A. DeSimone, M.C. Chuah, and O. Yue. Throughput performance of transport-layer protocols over wireless LANs. In *Proceedings of IEEE GLOBECOM*, 1993.
- [7] S. Floyd and T. Henderson. The newreno modification to tcp's fast recovery algorithm. Internet Request for Comments 2582, April 1999.
- [8] S. Keshav. *Congestion Control in Computer Networks*. PhD thesis, UC Berkeley, September 1991.
- [9] T. V. Lakshman and U. Madhow. The performance of TCP/IP for networks with high bandwidth-delay products and random loss. *IEEE/ACM Trans. Networking*, June 1997.
- [10] M. Mathis, J. Mahdavi, S. Floyd, and A. Romanow. TCP selective acknowledgement options. Internet Request for Comments 2018, October 1996.
- [11] S. Shenker. Some conjectures on the behavior of acknowledgment-based transmission control of random access communication channels. In *Proc. of ACM Sigmetrics*, 1987.

- [12] Prasun Sinha, Narayanan Venkitaraman, Raghupathy Sivakumar, and Vaduvur Bharghavan. WTCP: A reliable transport protocol for wireless wide-area networks. TIMELY Group Research Report, January 1999.
- [13] R. Yavatkar and N. Bhagawat. Improving end-to-end performance of tcp over mobile internetworks. In *IEEE Workshop on Mobile Computing Systems and Applications*, 1994.