Network Working Group                                M. Sridharan
Internet Draft                                          Microsoft
Intended status: Experimental July 18, 2007                K. Tan
Expires: January 2008                          Microsoft Research
                                                        D. Bansal
                                                        D. Thaler
                                                        Microsoft

Compound TCP: A New TCP Congestion Control for High-Speed and Long
                        Distance Networks


                 draft-sridharan-tcpm-ctcp-00.txt


Status of this Memo

  Internet-Drafts are working documents of the Internet Engineering
  Task Force (IETF), its areas, and its working groups.  Note that
  other groups may also distribute working documents as Internet-
  Drafts.

  Internet-Drafts are draft documents valid for a maximum of six months
  and may be updated, replaced, or obsoleted by other documents at any
  time.  It is inappropriate to use Internet-Drafts as reference
  material or to cite them other than as "work in progress."

  The list of current Internet-Drafts can be accessed at
  http://www.ietf.org/ietf/1id-abstracts.txt.

  The list of Internet-Draft Shadow Directories can be accessed at
  http://www.ietf.org/shadow.html.

   This Internet-Draft will expire on January 18, 2008.

Copyright Notice

Abstract

This document proposes Compound TCP (CTCP), a modification to TCP's congestion control mechanism for use with TCP connections with large congestion windows. The key idea behind CTCP is to add a scalable delay-based component to the standard TCP's loss-based congestion control. The sending rate of CTCP is controlled by both loss and delay components. The delay-based component has a scalable window increasing rule that not only efficiently uses the link capacity, but on sensing queue build up, gracefully reduces the sending rate. We have implemented CTCP on Microsoft's Windows and we have done extensive testing on production links and in Windows Beta deployments. We also engaged with Stanford Linear Accelerator Center to evaluate the properties of CTCP. The results so far are very encouraging. This document describes the Compound TCP algorithm in detail, and solicits experimentation and feedback from the wider community. In this document, we collectively refer to any TCP congestion control algorithm that employs a linear increase function for congestion control, including TCP Reno and all its variants as Standard TCP.

Table of Contents

## 1. Introduction

This document proposes Compound TCP, a modification to TCP's congestion
control mechanism for fast, long-distance networks. The standard TCP
congestion avoidance algorithm employs an additive increase and
multiplicative decrease (AIMD) scheme, which employs a conservative
linear growth function for increasing the congestion window and
multiplicative decrease function on encountering a loss. For a high-
speed and long delay network, it will take standard TCP an unreasonably
long time to recover the sending rate after a single loss event
[RFC2581, RFC3649]. Moreover, it is well-known now that in a steady-
state environment, with a packet loss rate of p, the current standard
TCP's average congestion window is inversely proportional to the square
root of the packet loss rate [RFC2581,PADHYE]. Therefore, it requires
an extremely small packet loss rate to sustain a large window. As an
example, Floyd et al. [RFC3649], pointed out that under a 10Gbps link
with 100ms delay, it will roughly take one hour for a standard TCP flow
to fully utilize the link capacity, if no packet is lost or corrupted.
This one hour error free transmission requires a packet loss rate
around $10^{-11}$ with 1500-byte size packets (one packet loss over
2,600,000,000 packet transmission!), which is not practical in today's
networks.

There are several proposals to address this fundamental limitation of
TCP. One straightforward way to overcome this limitation is to modify
TCP control's increase/decrease rule in its congestion avoidance stage.
More specifically, in the absence of packet loss, the sender increases
congestion window more quickly and decreases it more gently upon a
packet loss. In a mixed network environment, the aggressive behavior of
such approaches may severely degrade the performance of regular TCP
flows whenever the network path is already highly utilized. When an
aggressive high-speed variant flow traverses the bottleneck link with
other standard TCP flows, it may increase its own share of bandwidth by
reducing the throughput of other competing TCP flows. As a result the
aggressive variants will cause much more self-induced packet losses on
bottleneck links, and push back the throughput of the regular TCP
flows.

Then there is the class of high-speed protocols which use variances in
RTT as a congestion indicator (e.g., [AFRICA,FAST]). The delay-based
approaches are more-or-less derived from the seminal work of TCP-Vegas
[VEGAS]. An increase in RTT is considered an early indicator of
congestion, and the sending rate is cut in half to avoid buffer
overflow. The problem in this approach comes when delay-based and loss-
based flows share the same bottleneck link. While the delay-based flows
respond to increases in RTT by cutting its sending rate, the loss-based
flows continue to increase their sending rate. As a result a delay-
based flow obtains far less bandwidth than its fair share. This
weakness is hard to remedy for purely delay based approaches.

The design of Compound TCP is to satisfy to efficiency requirement and
TCP friendliness requirement simultaneously. The key idea is that if
the link is under-utilized, the high-speed protocol should be
aggressive and increase the sending rate quickly. However, once the
link is fully utilized, being aggressive will not only adversely affect
standard TCP flows but will also cause instability. As noted above,
delay-based approaches already have a nice property of adjusting its
aggressiveness based on the link utilization, which is observed by the
end-systems as an increase in RTT. CTCP incorporates a scalable delay-
based component to the standard TCP's congestion avoidance algorithm.
Using the delay component as an automatic tuning knob, CTCP is scalable
yet TCP friendly.

## [2](). Design Goals

The design of CTCP is motivated by the following requirements:

    o  Improve throughput by efficiently using the spare capacity in
       the network
    o  Good intra-protocol fairness when competing with flows that
       have different RTTs
    o  Should not impact the performance of standard TCP flows sharing
       the same bottleneck
    o  No additional feedback or support required from the network

CTCP can efficiently use the network resource and achieve high link
utilization. The aggressiveness can be controlled by adopting a rapid
increase rule in the delay-based component. We choose CTCP to have
similar aggressiveness as HighSpeed TCP [RFC3649]. Our design choice is
motivated by the fact that HSTCP has been tested to be aggressive
enough in real world networks and is now an experimental IETF RFC. We
also wanted an upper bound on the amount of unfairness to standard TCP
flows. However, as shown later, CTCP is able to maintain TCP
friendliness under high statistical multiplexing and also while
traversing poorly buffered links. CTCP has similar or in some cases,
even improved RTT fairness compared to standard TCP. As we will
demonstrate later this is due to the fact that the amount of backlogged
packets for a connection is independent of the RTT of the connection.
Even though CTCP does not require any feedback from the network, CTCP
works well in ECN capable environments. There is also no expectation on
the queuing algorithm deployed in the routers.

As is the case with most high-speed variants today, CTCP does not
modify slow-start. We agree to the belief that ramping-up faster than
slow-start without additional information from the network can be
harmful. Similar to HSTCP, to ensure TCP compatibility, CTCP's scalable
component uses the same response function as Standard TCP when the
current congestion window is at most Low_Window. CTCP sets Low_Window
to 38 MSS-sized segments, corresponding to a packet drop rate of 10^-3
for TCP.

## [3](). Compound TCP Control Law

CTCP modifies Standard TCP's loss-based control law with a scalable
delay-based component. To do so, a new state variable is introduced in
current TCP Control Block (TCB), namely, dwnd (Delay Window), which
controls the delay-based component in CTCP. The conventional congestion
window, cwnd, remains untouched, which controls the loss-based
component in CTCP. Thus, the CTCP sending window now is controlled by

both cwnd and dwnd. Specifically, the TCP sending window (wnd) is now
calculated as follows:

   wnd = min(cwnd + dwnd, awnd),              (1)

where awnd is the advertised window from the receiver.

cwnd is updated in the same way as regular TCP in the congestion
avoidance phase, i.e., cwnd is increased by 1 MSS every RTT and halved
when a packet loss is encountered. The update to dwnd will be explained
in detail later in the section. The combined window for CTCP from (1)
above allows up to (cwnd + dwnd) packets in one RTT. Therefore, the
increment of cwnd on the arrival of an ACK is modified accordingly:

   cwnd = cwnd + 1/(cwnd+dwnd)                (2)

As stated above, CTCP retains the same behavior during slow start. When
a connection starts up dwnd is initialized to zero while the connection
is in slow start phase. Thus the delay component is effective when the
connection enters congestion avoidance. The delay-based algorithm has
the following properties. It uses a scalable increase rule when it
infers that the network is under-utilized. It also reduces the sending
rate when it sense incipient congestion. By reducing its sending rate,
the delay-based component yields to competing TCP flows and ensures TCP
fairness. It reacts to packet losses by reducing its sending rate,
which is necessary to avoid congestion collapse. Our control law for
the delay-based component is derived from TCP Vegas. A state variable,
called basertt tracks the minimum round trip delay seen by a packet
over the network path. When a connection is started, basertt is updated
to be the minimum RTT observed during the 3-way handshake. The CTCP
sender also maintains a smoothed RTT srtt, updated as specified in
[RFC2988]. Then, the number of backlogged packets of the connection can
be estimated using,

   expected (throughput) = wnd/basertt
   actual (throughput) = wnd/srtt
   diff = (expected - actual) * basertt

The expected throughput gives the estimation of throughput CTCP gets if
it does not overrun the network path. The actual throughput stands for
the throughput CTCP really gets. Using this we can calculate the amount

of data backlogged in the bottleneck queue (diff). Congestion is detected by comparing diff to a threshold gamma. If diff < gamma, the network path is assumed to be under-utilized; otherwise the network path is assumed to be congested and CTCP should gracefully reduce its window.

It is to be noted that a connection should have at least gamma packets backlogged in the bottleneck queue to be able to detect incipient congestion. This motivates the need for gamma to be small since the implication is that even when the bottleneck buffer size is small, CTCP will react early enough to ensure TCP fairness. On the other hand if gamma is too small compared to the queue size, CTCP will falsely detect congestion and will adversely affect the throughput. Choosing the appropriate value for gamma could be a problem because this parameter depends on both network configuration and the number of concurrent flows, which are generally unknown to the end-systems. We present an effective way to automatically estimate gamma later in later sections.

The increase law of the delay-based component should make CTCP more scalable in high-speed and long delay pipes. We choose a binomial function to increase the delay window [BAINF01]. More specifically, when no congestion is detected, CTCP window increases using the following function

$$dwnd(t+1) = dwnd(t) + alpha*dwnd(t)^k \quad (3)$$

When a packet loss occurs, the delay window is multiplicatively decreased,

$$dwnd(t+1) = dwnd(t)*(1-beta) \quad (4)$$

where alpha, beta and k are tunable to obtain the desirable scalability, smoothness and responsiveness. We assume that a loss is detected by three duplicate ACKs. As explained in the next section we have modeled the response function for CTCP to have comparable scalability to HighSpeed TCP. Since there is already a loss-based component in CTCP, the delay-based component needs to be designed to only fill the gap, and the overall CTCP should follows the behavior defined in (3) and (4). We now summarize the control law for CTCP's delay component as follows;

```
 dwnd(t+1) =
     dwnd(t) + alpha*dwnd(t)^k - 1,      if diff < gamma  (5)
     dwnd(t) - eta*diff,                 if diff >= gamma (6)
     dwnd(t)(1-beta) - cwnd/2,           on packet loss   (7)
```

where (5) shows that in the increase phase, dwnd only needs to increase
by (alpha*dwnd(t)^k - 1) packets, since the loss-based component cwnd
will also increase by 1 packet. When a packet loss occurs, dwnd is set
to the difference between the desired reduced window size and that can
be provided by cwnd. The rule in equation (6) is very important to
preserve good RTT and TCP fairness. Eta defines how rapidly the delay
component should reduce its window when congestion is detected. Note
that dwnd is never negative, so the CTCP window is lower bounded by its
loss based component, which is same as Standard TCP.

If a retransmission timeout occurs, dwnd should be reset to zero and
the delay-based component is disabled. It is because that after a
timeout, the TCP sender enters slow-start phase. After the CTCP sender
exits the slow-start recovery state and enters congestion avoidance,
dwnd control kicks in again.


## 4. Compound TCP Response Function

The TCP response function provides a relationship between TCP's average
congestion window w in MSS-sized segments as a function of the steady-
state packet drop rate p. To specify a modified response function for
CTCP, we use the analytical model in [CTCPI06] to derive a relationship
between w and p. Based on this model, the response function for CTCP
provides the following relationship between w and p,

$$w \sim .1/(p^{(1/2-k)})  \qquad (8)$$

As explained earlier we modeled the response function for CTCP to have
comparable scalability to HighSpeed TCP. The response function for
HighSpeed TCP is

$$w \sim .1/p^{0.835}  \qquad (9)$$

Comparing (8) and (9) we get k to be around 0.8. Since it's difficult
to implement an arbitrary power we choose k = 0.75 which can be
implemented using a fast integer algorithm for square root. Based on
extensive experimentation, we choose alpha = 1/8 and beta = 1/2.
Substituting the above values for alpha, beta and k in (8) we get the

following response function for CTCP,

    w = 0.255/p^0.8        (10)

The response function for CTCP is compared with HSTCP and is
illustrated in Table 1 below.

|  | CTCP | HSTCP |
| --- | --- | --- |
| Packet Drop Rate P | Congestion Window W | Congestion Window W |
| ------------------ | ------------------- | ------------------- |
| 10^-3 | 64 | 38 |
| 10^-4 | 404 | 263 |
| 10^-5 | 2552 | 1795 |
| 10^-6 | 16107 | 12279 |
| 10^-7 | 101630 | 83981 |
| 10^-8 | 641245 | 574356 |
| 10^-9 | 4045987 | 3928088 |
| 10^-10 | 25528453 | 26864653 |

    Table 1: TCP Response function for CTCP & HSTCP

The values in Table 1 illustrate that our choice of parameters makes
CTCP slightly more aggressive than HSTCP in moderate and low packet
loss rates but approaches HSTCP for larger windows. The reason we
choose to do this is because unlike HighSpeed TCP, CTCP's delay control
is capable of scaling back on detecting incipient congestion. As a
result we expect CTCP to be more TCP friendly than HighSpeed TCP. We
show that this is in fact the case even under low buffering conditions
in the presence of high statistical multiplexing. The fairness
considerations and choice of gamma are detailed in later sections.

## 5. Automatic Selection of Gamma

To effectively detect early congestions, CTCP requires estimating the
backlogged packets at bottleneck queue and compares this estimate to a
pre-defined threshold gamma. However, setting this threshold gamma is
particular difficult for CTCP (and to many other similar delay-based
approaches), because gamma largely depends on the network configuration
and the number of concurrent flows that compete for the same bottleneck
link, which are, unfortunately, unknown to end-systems. Based on
experimentation over varying conditions we selected gamma to be 30
packets. This value provided a pretty good tradeoff between TCP
fairness and throughput. However a fixed gamma can still result in poor
TCP friendliness over under-buffered network links. One naive solution

is to choose a very small value for gamma, however this can falsely
detect congestion and adversely affect throughput. To address this
problem we use a method called tuning-by-emulation to dynamically
adjust gamma. The basic idea of our proposal is to estimate the
backlogged packets of a Standard TCP flow along the same path by
emulating the behavior of a Standard TCP flow in runtime. Based on
this, gamma is set so as to ensure good TCP-friendliness. CTCP can then
automatically adapt to different network configurations (i.e., buffer
provisioning) and also concurrent competing flows.

Our analytical model on CTCP shows that gamma should at least be less
than B/m+l to ensure the effectiveness of incipient congestion
detection, where m and l present the flow number of concurrent Standard
TCP flows and CTCP flows that are competing for the same bottleneck
link [CTCPI06,CTCPP06,CTCPT]. Generally, both B and (m+l) are unknown
to end-systems. It is very difficult to estimate these values from end-
systems in real-time, especially the number of flows, which can vary
significantly over time. Fortunately there is a way to directly
estimate the ratio B/m+l, even though the individual variables B or
(m+l) are hard to estimate. Let's first assume there are (m+l) regular
TCP flows in the network. These (m+l) flows should be able to fairly
share the bottleneck capacity in steady state. Therefore, they should
also get roughly equal share of the buffers at the bottleneck, which
should equal to B/m+l. For such a Standard TCP flow, although it does
not know either B or (m+l), it can still infer B/m+l easily by
estimating its backlogged packets, which is a rather mature technique
widely used in many delay-based protocols.  This brings us to the core
idea of CTCP's algorithm; CTCP lets the sender emulate the congestion
window of a Standard TCP flow. Using this emulated window, we can
estimate the buffer occupancy (Q) for a Standard TCP flow. Q can be
regarded as a conservative estimate of B/m+l assuming that the high
speed flow is more aggressive than Standard TCP. By choosing gamma <=
Q, we can ensure TCP fairness.

The implementation is actually trivial. This is because CTCP already
emulates Standard TCP as the loss-based component. We can simply
estimate the buffer occupancy of a competing Standard TCP flow from
state which CTCP already maintains. We choose an initial gamma = 30 and
Q is calculated as follows,

```
expected_reno (throughput) = cwnd/basertt
actual_reno (throughput) = cwnd/srtt
diff_reno = (expected - actual) * basertt
```

The difference between diff_reno and diff is simply that diff_reno is
computed only using the loss based component cwnd. Since Standard TCP
reaches its maximum buffer occupancy just before a loss, CTCP uses the
diff_reno value computed in the earlier round to calculate the gamma
for the next round. Whenever a loss happens, gamma is chosen to be less
than diff_reno and the sample values of gamma are updated using a
standard exponentially weighted moving average. The pseudocode to
calculate gamma is shown below. Here a round tracks every window worth
of data. We will provide more details on how to maintain a round in
Section 7.

```
  Initialization:
    diff_reno = invalid;
     Gamma = 30;

  End-of-Round:

     expected_reno = cwnd / baseRTT;
     actual_reno = cwnd / RTT;
     diff_reno = (Expected_reno-Actual_reno)*baseRTT;

  On-Packet-Loss:

  If diff_reno is valid then
     g_sample = 3/4*Diff_reno;
     gamma = gamma*(1-lamda)+ lamda*g_sample;
     if (gamma < gamma_low)
       gamma=gamma_low;
     else if (gamma > gamma_high)
       gamma=gamma_high;
     fi
     diff_reno = invalid;
  fi
```

The recommended values for gamma_low and gamma_high are 5 and 30
respectively. diff_reno is set to invalid to prevent using stale
diff_reno data when there are consecutive losses between which no
samples were taken.

## 6. Implementation Issues

The first challenge is to design a mechanism that can precisely track the changes in round trip time with minimal overhead, and can scale well to support many concurrent TCP connections. Naively taking RTT sample for every packet will obviously be an over-kill for both CPU and system memory, especially for high-speed and long distance networks where the congestion window can be very large. Therefore, CTCP needs to limit the number of samples taken, but without compromising on accuracy. In our implementation, we only take up to M sample per window of data. M is chosen to scale with the round trip delay and window size.

In order to further improve the efficiency in memory usage, we have developed a memory allocation mechanism to dynamically allocate sample buffers from a kernel fixed-size per-processor pool. The size should be chosen as a function of the available system memory. As the window size increases, M can be updated so that the samples are uniformly distributed over the window. As M gets updated more memory blocks are allocated and linked to the existing sample buffers. If the sending rate changes either due to network conditions or due to application behavior, the sample blocks are reclaimed to the global memory pool. This dynamic buffer management ensures the scalability of our implementation, so that it can work well even in a busy server which could host tens of thousands of TCP connections simultaneously. Note that it may also require high-resolution timer to time RTT samples.

The rest of the implementation is rather straightforward. We add two new state variables into the standard TCP Control Block, namely dwnd and basertt. The basertt is a value that tracks the minimum RTT sample measured seen so far and it is used as an estimation of the transmission delay of a single packet. Basertt is usually cleared if a retransmission timeout is hit. It is a good idea to re-measure the basertt incase the network conditions have changed. Following the common practice of high-speed protocols, CTCP reverts to standard TCP behavior when the window is small. Delay-based component only kicks in when cwnd is larger than some threshold, currently set to 38 packets assuming 1500 byte MTU. dwnd is updated at the end of each round. Note that no RTT sampling and dwnd update happens during the loss recovery phase. It is because the retransmission during the loss recovery phase may result in inaccurate RTT samples and can adversely affect the

delay-based control.

## [7]. Deployment Issues

There are several variations of TCP proposed for high speed and long
delay networks. We do not claim Compound TCP to be the best nor the
most optimal algorithm. However, based on our extensive testing via
simulations, experimentation including those on production links as
well as beta deployments of a reasonable scale, we believe that
Compound TCP satisfies the design considerations outlined before in
this document. It effectively uses spare bandwidth in high speed
networks, achieves good intra-protocol fairness even in the presence of
differing RTTs and does not adversely impact standard TCP. Further,
Compound TCP does not require any changes or any new feedback from the
network and is deployable over the current Internet in an incremental
fashion. It inter-operates with Standard TCP and requires support only
one the send side of a TCP connection for it to be used.
We also note that similar to High Speed TCP, in environments typical of
much of the current Internet, Compound TCP behaves exactly like
Standard TCP. This it does by ensuring that is follows standard TCP
algorithm without any modification any time congestion window is less
than 38 packets. Only when congestion window is greater than 38
packets, does the delay based component of Compound TCP gets invoked.
Thus, for example for a connection with RTT of 100ms, end to end
bandwidth must be greater than 4.8Mbps for CTCP algorithm to have any
difference in its response to network conditions than a standard TCP.

Further, we do not believe that the deployment of Compound TCP would
block the possible deployment of alternate experimental congestion
control algorithms such as Fast TCP [FAST] or CUBIC [CUBIC]. In
particular, Compound TCP s response has a fallback to loss based
function that has characteristics very similar to HS-TCP or N parallel
TCP connections.

## [8].   Security Considerations

This proposal makes no changes to the underlying security of the TCP
protocol.

## [9].   IANA Considerations

There are no IANA considerations regarding this proposal.

## 10. Conclusions

This document proposes a novel congestion control algorithm for TCP for high speed and long delay networks. By introducing a delay based component in addition to a standard TCP based loss component, Compound TCP is able to detect and effectively use spare bandwidth that may be available on a high speed and long delay network. Further, delay based component detects onset of congestion early and gracefully reduces sending rate. The loss based component, on the other hand, ensures there is effective response to losses in network while in the absence of losses, keeps the throughput of CTCP lower bounded by TCP Reno. Thus, CTCP is not timid, nor induces more self induced packet loss than a single standard TCP flow. Thus Compound TCP is efficient in consuming available bandwidth while being friendly to standard TCP. Further, the delay component does not have any RTT bias thereby reducing the RTT bias of the Compound TCP vis-a-vis standard TCP.

Compound TCP has been implemented as an optional component in Microsoft Windows Vista Operating System. It has been tested and experimented through broad Windows Vista beta deployments where it has been verified to meet its objectives without causing any adverse impact. SLAC has also evaluated Compound TCP on production links. Based on testing and evaluation done so far, we believe Compound TCP is safe to deploy on the current Internet. We welcome additional analysis, testing and evaluation of Compound TCP by Internet community at large and continue to do additional testing ourselves.

## 11. Acknowledgments

The authors would like to thank Jingmin Song for all his efforts in evaluating the algorithm on the test beds. We are thankful to Yee-ting Lee and Les Cottrell for testing and evaluation of Compound TCP on Internet2 links [SLAC]. We would like to thank Sanjay Kaniyar for his insightful comments and for driving this project in Microsoft. We are also thankful to the Microsft.com data center staff who helped us evaluate Compound TCP on their production links. In addition, several folks from the Internet research community who attended the High-Speed TCP Summit at Microsoft [MSWRK] have provided valuable feedback on Compound TCP. Finally, we are thankful to the Windows Vista program beta participants who helped us test and evaluate CTCP.

12.  References

12.1. Normative References

   [RFC2581] Allman, M., Paxson, V. and W. Stevens, "TCP Congestion
             Control", RFC 2581, April 1999.

12.2. Informative References

   [AFRICA]  R. King, R. Baraniuk and R. riedi, "TCP-Africa: An
             Adaptive and Fair Rapid Increase Rule for Scalable
             TCP", In Proc. INFOCOM 2005.

   [BAINF01] D. Bansal and H. Balakrishnan, "Binomial Congestion
             Control Algorithms", Proc INFOCOM 2001.

   [CTCPI06] K. Tan, Jingmin Song, Qian Zhang, Murari Sridharan, "A
             Compound TCP Approach for High-speed and Long Distance
             Networks", in IEEE Infocom, April 2006, Barcelona,
             Spain.

   [CTCPP06] K. Tan, J. Song, Q. Zhang, and M. Sridharan, "Compound
             TCP: A Scalable and TCP-friendly Congestion Control for
             High-speed Networks", in 4th International workshop on
             Protocols for Fast Long-Distance Networks (PFLDNet),
             2006, Nara, Japan.

   [CTCPT]   K. Tan, J. Song, M. Sridharan, and C.Y. Ho, "CTCP:
             Improving TCP-Friendliness Over Low-Buffered Network
             Links", Microsoft Technical Report.

   [CUBIC]   I. Rhee, L. Xu and S. Ha, "CUBIC for fast long distance
             networks", Internet Draft, Expires Aug 31, 2007, draft-
             rhee-tcp-cubic-00.txt

   [FAST]    C. Jin, D. Wei, S. Low, "FAST TCP: Motivation,
             Architecture, Algorithms, Performance", in IEEE Infocom
             2004.

   [MSWRK]   Microsoft High-Speed TCP Summit,
             http://research.microsoft.com/events/TCPSummit/

   [PADHYE]  J. Padhya, V. Firoiu, D. Towsley and J. Kurose, "Modeling
             TCP Throughput: A Simple Model and its Empirical
             Validation", in Proc. ACM SIGCOMM 1998.

   [RFC2988] V. Paxson and M. Allman, "Computing TCP's Retransmission
             Timer", RFC 2988, November 2000.

   [RFC3649] S. Floyd, "HighSpeed TCP for Large Congestion Windows",
             RFC 3649, Dec 2003.

   [SLAC]    Yee-Ting Li, "Evaluation of TCP Congestion Control
             Algorithms on the Windows Vista Platform", SLAC-TN-06-
             005,
             http://www.slac.stanford.edu/pubs/slactns/tn04/slac-tn-
             06-005.pdf

   [VEGAS]   L. Brakmo, S. O'Malley, and L. Peterson, "TCP Vegas: New
             techniques for congestion detection and avoidance", in
             Proc. ACM SIGCOMM, 1994.

Authors' Addresses

   Murari Sridharan
   Microsoft Corporation
   1 Microsoft Way, Redmond 98052


   Email: muraris@microsoft.com


   Kun Tan
   Microsoft Research
   5/F, Beijing Sigma Center
   No.49, Zhichun Road, Hai Dian District
   Beijing China 100080

   Email: kuntan@microsoft.com


   Deepak Bansal
   Microsoft Corporation
   1 Microsoft Way, Redmond 98052

   Email: dbansal@microsoft.com


   Dave Thaler
   Microsoft Corporation
   1 Microsoft Way, Redmond 98052

   Email: dthaler@microsoft.com

Intellectual Property Statement

Disclaimer of Validity

Copyright Statement

Acknowledgment