

draft-leith-tcp-htcp-06
Internet-Draft
Intended status: Experimental
Expires: October 9, 2008

D. Leith
Hamilton Institute
April 7, 2008

H-TCP: TCP Congestion Control for High Bandwidth-Delay Product Paths

Status of this Memo

By submitting this Internet-Draft, each author represents that any applicable patent or other IPR claims of which he or she is aware have been or will be disclosed, and any of which he or she becomes aware will be disclosed, in accordance with [Section 6 of BCP 79](#).

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/lid-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 October 9, 2008.

Internet-Draft

H-TCP

April 2008

Abstract

This document describes a number of changes to the TCP congestion control algorithm to improve performance in high bandwidth-delay product paths. We focus on changes to the congestion avoidance mode, rather than slow-start.

Table of Contents

1.	Conventions	3
2.	Introduction	4
3.	Additive Increase for High Bandwidth-Delay Product Paths	6
4.	Impact of Changes on Performance	8
4.1.	RTT unfairness	8
4.2.	Friendliness	8
4.3.	Responsiveness	8
4.4.	Efficiency	8
5.	Optional RTT Scaling	10
6.	Security Considerations	11
7.	Acknowledgements	12
8.	Changelog	13
9.	Informative References	14
	Author's Address	15
	Intellectual Property and Copyright Statements	16

1. Conventions

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [RFC 2119](#).

2. Introduction

This document describes a number of changes to the TCP congestion control algorithm to improve performance in high bandwidth-delay product paths.

The current TCP congestion control algorithm is known to perform poorly on paths where the TCP congestion window becomes large. [[Kelly02](#), [Flo03](#), [FAST04](#)]. Following congestion, the congestion window is halved and only increases at a rate of 1 packet per RTT. As a result flows can take an unacceptably long time to recover their window size after a congestion event.

A direct solution is to make the time between congestion events smaller. This can be achieved by, for example, adjusting the AIMD additive increase rate to be greater for flows with larger congestion window. Backward compatibility with legacy TCP can be ensured through the inclusion of a separate mode of operation that behaves as legacy TCP in the appropriate circumstances.

The logic that orchestrates switching between the legacy and more aggressive modes of operation can clearly be designed several ways. One approach is to make the AIMD increase parameter, which we denote here by α , a function of the flow congestion window. That is, α is increased as congestion window increases thereby resulting in an additive increase algorithm that directly scales with congestion window. This is precisely the approach adopted in the High-Speed TCP [[Flo03](#)] proposal. In addition to adjusting the AIMD increase parameter α as a function of congestion window, this

proposal also increases the multiplicative decrease factor beta to further increase the aggressiveness of a flow. (Note. On multiplicative decrease, the congestion window cwnd is updated to beta x cwnd. We use this definition of the backoff factor beta throughout this document).

While such modifications might appear straightforward, it has been shown [[Sho04](#), [Yi05](#)] that they often negatively impact the behaviour of networks of TCP flows. High-speed TCP [[Flo03](#)], BIC-TCP [[BIC04](#)] and Cubic can exhibit slow convergence following network disturbances such as the start-up of new flows; Scalable-TCP [[Kelly02](#)] is a multiplicative-increase multiplicative-decrease strategy and as such it is known that it may fail to converge to fairness in drop-tail networks [[Jain89](#)].

Scope

Our focus in this document is on the behaviour of long-lived flows

Leith

Expires October 9, 2008

[Page 4]

Internet-Draft

H-TCP

April 2008

and so we do not consider changes to slow-start. We also seek to make the smallest possible changes to the existing TCP congestion control algorithm, and so confine consideration to the AIMD packet-loss based paradigm. Use of jumbo packets is viewed as complementary to the changes proposed here.

3. Additive Increase for High Bandwidth-Delay Product Paths

It is known that modifying the AIMD backoff factor can have a significant impact on network responsiveness, and this is discussed in more detail elsewhere [[Sho04](#), [Sho05](#)]. In this document we confine attention to modifications to the AIMD increase rate with the aim of improving performance in high bandwidth-delay product paths. We begin with the observation that making the AIMD increase rate an increasing function of flow cwnd (as is done in the HS-TCP, BIC, Cubic etc algorithms) means that flows with smaller cwnd are placed at a disadvantage to flows with larger cwnd when competing for bandwidth. This is a primary source of unfairness and slow convergence. We therefore take an alternative approach. Noting that it is the increase in congestion epoch duration with bandwidth-delay product that is the source of many issues, we make the AIMD increase

rate purely a function of the elapsed time since the last congestion event. This allows us to increase the aggressiveness of the AIMD increase as the congestion epoch duration increases (so improving performance in high bandwidth-delay product paths) while avoiding placing flows with small cwnd at a consistent disadvantage.

[RFC2591](#) specifies that during congestion avoidance, cwnd is incremented by 1 full-sized segment per round-trip time (RTT). We modify this behaviour to increase cwnd by alpha segments per RTT, where alpha is calculated as follows.

```
if Delta <= Delta_L
    alpha = 1
else
    alpha = f_alpha(Delta)
```

where Delta is the time in seconds that has elapsed since the last congestion event experienced by a flow and Delta_L is the threshold for switching from standard/legacy operation to the new increase function. Delta_L MUST be at least 1 second, although larger values MAY be used. The increase function f_alpha is selected such that the duration of the congestion epochs remains reasonably small as the bandwidth-delay product on a path increases. Below, we discuss a choice of increase function that yields convergence times that seem reasonable. However, the precise responsiveness requirement in future networks is currently not well defined and so the specific choice of increase function may change.

Use of the following increase function is RECOMMENDED:

$$f_alpha(Delta) = 1 + 10(Delta-Delta_L)+0.5(Delta-Delta_L)^2 \quad (1)$$

This choice yields the congestion epoch duration for a single flow,

as a function of congestion window size, shown in Table 1.

-----	-----
Congestion window (packets)	Congestion epoch duration (s)
-----	-----
100	1.1

1000	3.1
2000	4.3
5000	6.6
10000	9.2
20000	12.8
50000	19.4

Table 1 - Congestion epoch duration vs congestion window size for an RTT of 100ms

[4.](#) Impact of Changes on Performance

[4.1.](#) RTT unfairness

The level of unfairness between flows with different RTT's is similar to that with the standard TCP algorithm. This behaviour is confirmed in experimental and simulation tests [[HTCP04](#), [Yi05](#)].

[4.2.](#) Friendliness

The mean AIMD increase parameter is shown in Table 2 for a range of bandwidth-delay products. This an indication of the number of standard TCP flows (neglecting statistical multiplexing of backoffs) whose aggregate would be equivalent to a flow using increase function (1). That is, an indication of friendliness and also of the packet drop overhead associated with the AIMD probing action.

Congestion window (packets)	Effective number of standard TCP flows		
	10ms RTT	100ms RTT	250ms RTT
10	1	1	1
100	1	2	5
1000	3	12	22
2000	4	19	32
5000	8	33	55
10000	12	49	82
20000	19	72	123
50000	32	122	208

Table 2 - Mean increase parameter (packets/RTT) vs congestion window size

[4.3.](#) Responsiveness

Responsiveness is qualitatively similar to that of the current AIMD congestion control algorithm, i.e. the convergence time of TCP flows using an AIMD backoff factor of 0.5 is approximately 4 congestion epochs, although the congestion epoch duration is significantly shorter on high bandwidth-delay product paths (see Table 1).

[4.4.](#) Efficiency

Link utilisation depends on queue provisioning in a similar manner to the current TCP congestion control algorithm. That is, for a single flow (or multiple synchronised flows) 100% link utilisation requires that the queue be sized as the bandwidth-delay product. Simulation

and experimental tests indicate that statistical multiplexing between unsynchronised flows yields similar efficiency gains to standard TCP.

5. Optional RTT Scaling

We note that the parameter α determines the AIMD increase rate in packets per RTT. Hence, flows with the same RTT have the same increase rate in packets per second, but flows with different RTTs have different increase rate in packets per second. It is this that primarily leads to unfairness between flows with different RTTs. Removing RTT unfairness is not one of our objectives here. However, we note that an AIMD flow generates roughly α packet drops per RTT as a result of its probing action. Hence, flows with short RTT are more aggressive than flows with long RTT in the sense that they generate more packet drops over intervals of time measured in seconds. We can reduce the aggressiveness of short RTT flows by scaling the increase parameter α with RTT. This need not compromise the responsiveness of TCP flows. As noted in [Sh04, Sh05, HTCP04], the convergence time of TCP flows using an AIMD backoff factor of 0.5 is approximately 4 congestion epochs. Scaling α by RTT leads to scaling of the congestion epoch duration to become effectively the same for both short and long RTT flows. The convergence time is therefore also scaled to be effectively the same for both short and long RTT flows.

Such RTT scaling MAY be implemented by modifying the increase rule to

```
if Delta <= Delta_L
    alpha = 1
else
    alpha = K x f_alpha(Delta)
```

where $K = RTT/RTT_{ref}$. Note that RTT scaling is not applied in low-speed conditions in order to maintain backward compatibility with legacy TCP flows (ensuring adequate backward compatibility presented a major difficulty in previous studies on the use of RTT scaling). Note also that the scaling is proportional to RTT rather than RTT^2 , as we do not seek to achieve throughput fairness here. RTT_{ref} is the reference RTT for which f_{α} is designed to ensure acceptable congestion epoch durations, with the recommended value being 100ms.

Leith

Expires October 9, 2008

[Page 10]

Internet-Draft

H-TCP

April 2008

[6.](#) Security Considerations

Security implications are not discussed in this document.

Leith

Expires October 9, 2008

[Page 11]

Internet-Draft

H-TCP

April 2008

[7.](#) Acknowledgements

This work was supported by Science Foundation Ireland grants 00/PI.1/C067 and 04/IN3/I460.

Leith

Expires October 9, 2008

[Page 12]

Internet-Draft

H-TCP

April 2008

[8.](#) Changelog

April 2008: Updated to use [RFC2119](#) terminology. Discussion streamlined.

9. Informative References

[Jain89] D.M. Chiu, R. Jain, Analysis of the increase and decrease algorithms for congestion avoidance in computer networks. Computer Networks and ISDN Systems, 1989.

[Flo03] S.Floyd, HighSpeed TCP for Large Congestion Windows . Sally Floyd. IETF [RFC 3649](#), Experimental, Dec 2003.

[FAST04] C. Jin, D.X. Wei, S.H. Low, FAST TCP: motivation, architecture, algorithms, performance. Proc IEEE INFOCOM 2004.

[Kelly02] T. Kelly, On engineering a stable and scalable TCP variant, Cambridge University Engineering Department Technical Report CUED/F-INFENG/TR.435, June 2002.

[HTCP04] D.J.Leith, R.N.Shorten, H-TCP Protocol for High-Speed Long-Distance Networks. Proc. 2nd Workshop on Protocols for Fast Long Distance Networks. Argonne, USA, 2004.
<http://www.hamilton.ie/net/htcp3.pdf>

[BIC04] L. Xu, K. Harfoush, I. Rhee, Binary Increase Congestion Control for Fast Long-Distance Networks. Proc. INFOCOM 2004.

[Sho04] R.N.Shorten, D.J.Leith,J.Foy, R.Kilduff, Analysis and design of congestion control in synchronised communication networks. Automatica, 2004. <http://www.hamilton.ie/net/synchronised.pdf>

[Sho05] R.N.Shorten, F. Wirth,F., D.J. Leith, A positive systems model of TCP-like congestion control: Asymptotic results.
http://www.hamilton.ie/net/unsynchronised_final.pdf

[Yi05] Y.Li, D.J.Leith, R.N.Shorten, Experimental evaluation of TCP protocols of high-speed networks. <http://www.hamilton.ie/net/eval/>

[Cot05] R.L. Cottrell, S. Ansari, P. Khandpur, R. Gupta, R. Hughes-Jones, M. Chen, L. MacIntosh, F. Leers, Characterization and Evaluation of TCP and UDP-Based Transport On Real Networks. . Proc. 3rd Workshop on Protocols for Fast Long-distance Networks, Lyon, France, 2005.

[Hegde04] S. Hegde, D. Lapsley, B. Wyrowski, J. Lindheim, D.Wei, C. Jin, S. Low, H. Newman, FAST TCP in High Speed Networks: An Experimental Study. Proc. GridNets, San Jose, 2004.

NUI Maynooth
Maynooth, Co. Kildare
Ireland

Email: doug.leith@nuim.ie

Full Copyright Statement

Copyright (C) The IETF Trust (2008).

This document is subject to the rights, licenses and restrictions contained in [BCP 78](#), and except as set forth therein, the authors retain all their rights.

This document and the information contained herein are provided on an "AS IS" basis and THE CONTRIBUTOR, THE ORGANIZATION HE/SHE REPRESENTS OR IS SPONSORED BY (IF ANY), THE INTERNET SOCIETY, THE IETF TRUST AND THE INTERNET ENGINEERING TASK FORCE DISCLAIM ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

Intellectual Property

The IETF takes no position regarding the validity or scope of any Intellectual Property Rights or other rights that might be claimed to pertain to the implementation or use of the technology described in this document or the extent to which any license under such rights might or might not be available; nor does it represent that it has made any independent effort to identify any such rights. Information on the procedures with respect to rights in RFC documents can be found in [BCP 78](#) and [BCP 79](#).

Copies of IPR disclosures made to the IETF Secretariat and any assurances of licenses to be made available, or the result of an attempt made to obtain a general license or permission for the use of such proprietary rights by implementers or users of this specification can be obtained from the IETF on-line IPR repository at <http://www.ietf.org/ipr>.

The IETF invites any interested party to bring to its attention any copyrights, patents or patent applications, or other proprietary rights that may cover technology that may be required to implement this standard. Please address the information to the IETF at ietf-ipr@ietf.org.

Leith

Expires October 9, 2008

[Page 16]