

Network Working Group
Internet-Draft
Intended status: Experimental
Expires: August 30, 2007

M. Scharf
University of Stuttgart
S. Floyd
ICIR
P. Sarolahti
Nokia Research Center
February 26, 2007

Avoiding Interactions of Quick-Start TCP and Flow Control
draft-scharf-tsvwg-quick-start-flow-control-00.txt

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 August 30, 2007.

Copyright Notice

Copyright (C) The IETF Trust (2007).

Abstract

This document describes methods to avoid interactions between the flow control of the Transmission Control Protocol (TCP) and the Quick-Start TCP extension. Quick-Start is an optional TCP congestion control mechanism that allows hosts to determine an allowed sending rate from feedback of routers along the path. With Quick-Start, data

transfers can start with a potentially large congestion window. In order to fully utilize the data rate determined by Quick-Start, the sending host must not be limited by the TCP flow control, i. e., the amount of free buffer space advertised by the receive window.

There are two potential interactions between Quick-Start and the TCP flow control: First, receivers might not provide sufficiently large buffer space after connection setup, or they may implement buffer allocation strategies that implicitly assume the slow-start behavior on the sender side. This document therefore provides guidelines for buffer allocation in hosts supporting the Quick-Start extension. Second, the TCP receive window scaling mechanism interferes with Quick-Start when being used in the initial three-way handshake connection setup. This document describes a simple solution to overcome this problem.

Table of Contents

1. Introduction	3
2. Requirements notation	4
3. Quick-Start TCP and receive buffer dimensioning	4
3.1. Receiver buffer allocation strategies	4
3.2. Recommendations for buffer dimensioning in presence of Quick-Start requests	4
4. Quick-Start TCP and receive window scaling	5
4.1. Receive window scaling	5
4.2. Problem within the three-way handshake	5
4.3. Possible remedy	6
4.4. Discussion and deployment considerations	8
5. Security Considerations	8
6. IANA considerations	9
7. Acknowledgments	9
8. References	9
8.1. Normative References	9
8.2. Informative References	10
Authors' Addresses	10
Intellectual Property and Copyright Statements	12

1. Introduction

Quick-Start is an experimental extension for the Transmission Control Protocol (TCP) [RFC0793] that allows to speed up best effort data transfers. The Quick-Start TCP extension is specified in [RFC4782]. With Quick-Start, TCP hosts can request permission from the routers along a network path to send at a higher rate than allowed by the default TCP congestion control, in particular after connection setup or longer idle periods. The explicit router feedback avoids the time-consuming capacity probing by the TCP slow-start and can significantly improve transfer times over paths with a high bandwidth-delay product [SAF07].

The usage of Quick-Start significantly changes the TCP behavior during connection setup. This is why special care is needed in order to prevent interactions between Quick-Start and other TCP mechanisms. Specifically, TCP flow control mechanisms have to be optimized for the usage of Quick-Start, in particular when the TCP connection spans a path with a large bandwidth-delay product (BDP). In such cases the sending window should have a large value in order to achieve good TCP performance (see [RFC2488],[RFC3481]).

Unlike the standard slow-start mechanism, the Quick-Start TCP extension allows the sender to use large congestion windows immediately after connection setup. The usage of such large windows raises two questions: First, what receiver buffer allocation strategies should be used in combination with Quick-Start? And second, how to appropriately signal these large windows? This document addresses these issues and shows that Quick-Start requires special mechanisms in both cases. The document thereby supplements the Quick-Start TCP specification [RFC4782], where flow control issues have not been addressed in detail.

The rest of this document is structured as follows: First, the question of receive buffer allocation in combination with Quick-Start is addressed and dimensioning guidelines are provided. Second, a modification of the receive window scaling mechanism [RFC1323] is specified, which is required to fully benefit from Quick-Start when the Quick-Start request is used in the initial <SYN> segment.

It should be noted that the effects and most methods discussed in this document are not specific to the Quick-Start TCP extension. They could also be used in combination with other proposals that cause a behavior more aggressive than standard TCP slow-start, for instance [LAJ+07].

2. Requirements notation

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 [RFC2119].

3. Quick-Start TCP and receive buffer dimensioning

3.1. Receiver buffer allocation strategies

The TCP sending window results from the minimum of the congestion window and the receive window (also called advertised receiver window) [RFC2581]. A small receive window prevents the TCP connection from fully utilizing paths with a larger bandwidth-delay product. As a consequence, on the one hand, a TCP receiver should advertise a receive window that is big enough to allow an efficient utilization of the connection path. On the other hand, hosts with a potentially high number of TCP connections need to optimize the buffer and memory usage to be able to serve a maximum possible number of TCP connections. Finding a fixed receive buffer size that is optimal between these two goals is difficult.

This is why many modern TCP implementations use an intelligent dynamic buffer management. There are different auto-tuning techniques and heuristics [Dun06] designed to prevent the receive window from limiting the data rate at the sender. An implementation using buffer size auto-tuning is described for instance in [SB05]. A common characteristic of most of these buffer allocation strategies is that they initially start with a rather small receive window. The more data arrives, the more buffer is allocated to the corresponding connection. This behavior is reasonable if the sender uses the standard slow-start algorithm and thus starts with a small congestion window anyway. However, when using Quick-Start, a large receive buffer may be required immediately after connection setup.

3.2. Recommendations for buffer dimensioning in presence of Quick-Start requests

When a host receives and approves a Quick-Start request, in particular during the connection setup, it SHOULD allocate a "reasonable" amount of buffer space so that a potential Quick-Start data transfer can start with a high sending window. If buffer size auto-tuning is used, it SHOULD be ensured that a sufficiently high initial receive window is announced. The handling of buffer space upon arrival of a Quick-Start request SHOULD be configurable by the corresponding application.

Determining an appropriate "reasonable" receive buffer size is not a trivial task and also depends on the available system resources. However, unlike standard TCP slow-start, the Quick-Start extension provides some additional information that could help to properly dimension the receive buffer. A reasonable buffer size would typically be a small multiple of the bandwidth-delay product of the path. An approximation of the available bandwidth can be directly obtained from the approved Quick-Start rate in the received request. If the round-trip time (RTT) to the Quick-Start originator is also known (e. g., if it has been cached from previous connections), a reasonable buffer size can be directly calculated as a small multiple of the BDP. In case that the round-trip time is not known, the buffer dimension could be done for a configurable "worst-case" RTT such as 500 ms.

4. Quick-Start TCP and receive window scaling

4.1. Receive window scaling

The TCP header specified in [RFC0793] uses a 16 bit field to report the receive window size to the sender. This effectively limits the sending window to 65 kB. To circumvent this problem, the "Window Scale" TCP extension [RFC1323] defines an implicit scale factor, which is used to multiply the window size value found in a TCP header to obtain a 32 bit window size. If enabled, the scale factor is announced during connection setup by the "Window Scale" TCP option in <SYN> and <SYN,ACK> segments.

In general, using receive window scaling is highly beneficial for TCP connections over path with a large bandwidth-delay product [RFC2488],[RFC3481]. Otherwise, the path capacity cannot fully be utilized by TCP. Quick-Start TCP can significantly speed up data transfers over such paths [RFC4782],[SAF07]. As a consequence, a host supporting Quick-Start SHOULD enable receive window scaling. If Quick-Start is used in the initial three-way handshake, the minimum required scaling factor can be obtained from the required receive buffer space, which can be approximated as described in the previous section.

4.2. Problem within the three-way handshake

A problem arises when the Quick-Start mechanism is used within the three-way handshake, and the Quick-Start request is added to the initial <SYN> segment: In this scenario, if the Quick-Start request is approved by the routers along the path, the receiver echoes back the Quick-Start response in the <SYN,ACK> segment. This process is illustrated in [RFC4782]. Upon reception of the <SYN,ACK> with the

Quick-Start response, the sender can set the congestion window to the determined value so that it can immediately start to send with the approved data rate.

However, [RFC1323] defines that the "Window field in a SYN (i.e., a <SYN> or <SYN,ACK>) segment itself is never scaled." This means that the maximum receive window that can be signaled to the sender in the <SYN,ACK> is 65 kB. As a consequence, the TCP flow control will prevent the TCP sender from having more than 65 kB of outstanding data, even if the receiver has much more free buffer, and the Quick-Start feedback allows a much larger congestion window.

This effect essentially limits the maximum amount of data sent by Quick-Start to 65 kB, when the sender sends the Quick-Start request in the initial <SYN> segment. Also, the congestion window after quitting the Quick-Start rate pacing phase is at most 65 kB, as the congestion window is set to the amount of outstanding data at this point. This is an undesirable restriction for the Quick-Start mechanism, even if 65 kB is still much more than the initial congestion window in slow-start that is allowed by [RFC3390].

This issue only occurs when Quick-Start is used in the three-way TCP connection setup procedure, and only in the direction of the client (connection originator) to the server. Still, this case is one of the planned usage scenarios for the Quick-Start TCP extension.

4.3. Possible remedy

The limitation imposed by the window scaling could be addressed in two different ways: First, one could deviate from [RFC1323] and use a scaled receive window in <SYN> and <SYN,ACK> segments, if they include Quick-Start options. This would avoid the problem sketched in the previous section, but it is not compliant with the TCP specification and the currently deployed TCP implementations.

This document describes a second, standard-compliant method: When a host receives a <SYN> segment with a Quick-Start option, it processes the option as described in [RFC4782]. Provided that the host has Quick-Start support enabled, the Quick-Start response is echoed back in the <SYN,ACK> segment. As explained, this segment cannot announce receive windows larger than 65 kB. If the receiver allocates a buffer space larger than 65 kB, an additional empty segment (without <SYN> flag) SHOULD be sent after the <SYN,ACK> segment, in order to announce the true receive window. The resulting message flow is depicted in Figure 1.

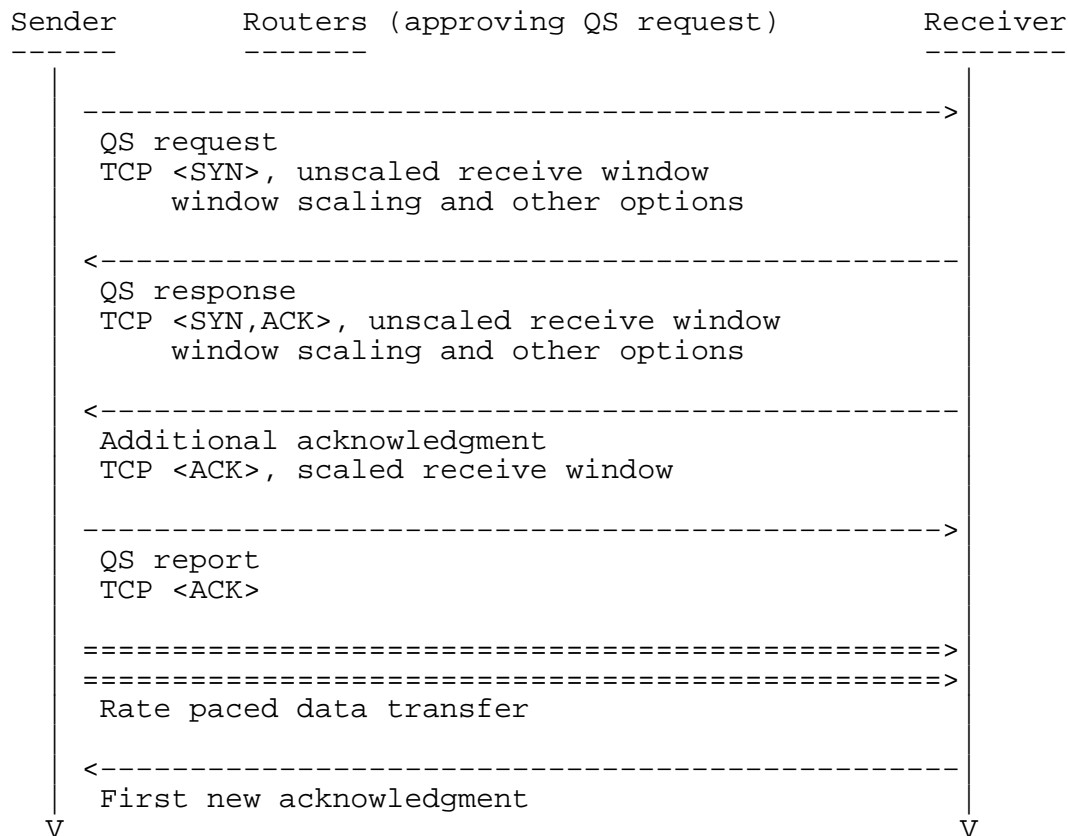


Figure 1: Message sequence chart of the proposed mechanism

After having received this additional acknowledgment, the sender is aware of the true available receive buffer. Provided that the Quick-Start request is approved on the path and that the receive window is sufficiently large, this allows the sender to send more than 65 kB during the Quick-Start rate pacing phase.

Note that there is some degree of freedom as to when to send the additional acknowledgment. It can be sent immediately after the <SYN,ACK> segment, but this is not required in all cases. It is sufficient if the sender receives this segment before reaching the limit of the unscaled receive window. As a consequence, receivers may decide to delay the sending of this segment for some small amount of time.

4.4. Discussion and deployment considerations

The method proposed in this document is compliant with the TCP specifications: Sending empty segments to increase the receive window is implicitly allowed by [RFC0793], and in [RFC2581] it is clearly stated that sending an acknowledgment is allowed to update the receive window. Implementing the method thus should require changes in the receiver TCP implementation only.

However, sending an empty acknowledgment shortly after a <SYN,ACK> segment is an atypical TCP communication event. The <SYN,ACK> and the additional segment could get reordered in the network. In this case, the sending host will typically ignore the additional segment, as it is still awaiting the <SYN,ACK>. Furthermore, middleboxes such as state-full firewalls might drop the additional acknowledgment. Even worse, this segment might also be dropped if a middlebox receives it earlier than the <ACK> segment from the sender. At this point in time, from the viewpoint of the middlebox, the bi-directional end-to-end TCP connection is not yet established. If the additional segment gets dropped, the sender gets informed about the unscaled receive window when the next new acknowledgment arrives, which may limit the benefit of Quick-Start. Delaying the additional acknowledgment for a short period of time could help to avoid such problems. Further investigation is needed to analyze whether such a delay is required.

A possible alternative to the message flow in Figure 1 would be to piggyback the Quick-Start response on the additional acknowledgment segment instead of the <SYN,ACK>. However, this approach has several drawbacks and is therefore not recommended: First, the Quick-Start response would be received later, which could cause additional delays. Second, the <SYN,ACK> is immediately acknowledged by the <ACK> segment. The Quick-Start rate report can thus be piggybacked on this <ACK>. In contrast, if the Quick-Start response is included in the additional acknowledgment, the Quick-Start report has to be piggybacked to a data segment, i. e., it depends on the availability of application data whether and when the Quick-Start report is sent.

It must be emphasized that the additional segment mandated by this document results in a certain network overhead. Given the fact that Quick-Start requests will be approved over under-utilized paths only, this overhead might not be a significant problem.

5. Security Considerations

Quick-Start TCP imposes a number of security challenges. Known security threats as well as counter-measures are discussed in the

section "Security Considerations" of [RFC4782]. Since this document describes extensions to Quick-Start TCP, the security issues identified in [RFC4782] apply here, too.

Sending an additional acknowledgment segment is an allowed behavior for a TCP connection endpoint and does not result in additional security threats. However, special care is needed when allocating large amounts of buffer space to newly established TCP connections, since this could create vulnerabilities to denial-of-service attacks. This issue may not be critical if Quick-Start is used in controlled environments only, as recommended by [RFC4782].

6. IANA considerations

This document has no actions for IANA.

7. Acknowledgments

The first author thanks Haiko Strotbek, Martin Koehn, Simon Hauger, and Christian Mueller for contributing to this document.

8. References

8.1. Normative References

- [RFC0793] Postel, J., "Transmission Control Protocol", STD 7, RFC 793, September 1981.
- [RFC1323] Jacobson, V., Braden, B., and D. Borman, "TCP Extensions for High Performance", RFC 1323, May 1992.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.
- [RFC2581] Allman, M., Paxson, V., and W. Stevens, "TCP Congestion Control", RFC 2581, April 1999.
- [RFC3390] Allman, M., Floyd, S., and C. Partridge, "Increasing TCP's Initial Window", RFC 3390, October 2002.
- [RFC4782] Floyd, S., Allman, M., Jain, A., and P. Sarolahti, "Quick-Start for TCP and IP", RFC 4782, January 2007.

8.2. Informative References

- [Dun06] Dunigan, T., "TCP auto-tuning zoo", available at <http://www.csm.ornl.gov/~dunigan/net100/auto.html>, February 2006.
- [LAJ+07] Liu, D., Allman, M., Jin, S., and L. Wang, "Congestion Control Without a Startup Phase", PFLDnet2007, Marina Del Rey, CA, USA, February 2007.
- [RFC2488] Allman, M., Glover, D., and L. Sanchez, "Enhancing TCP Over Satellite Channels using Standard Mechanisms", BCP 28, RFC 2488, January 1999.
- [RFC3481] Inamura, H., Montenegro, G., Ludwig, R., Gurtov, A., and F. Khafizov, "TCP over Second (2.5G) and Third (3G) Generation Wireless Networks", BCP 71, RFC 3481, February 2003.
- [SAF07] Sarolahti, P., Allman, M., and S. Floyd, "Determining an Appropriate Sending Rate Over an Underutilized Network Path", accepted for publication in Computer Networks, 2007.
- [SB05] Smith, M. and S. Bishop, "Flow Control in the Linux Network Stack", available at <http://www.cl.cam.ac.uk/~pes20/Netsem/linuxnet.pdf>, February 2005.

Authors' Addresses

Michael Scharf
University of Stuttgart
Pfaffenwaldring 47
D-70569 Stuttgart
Germany

Phone: +49 711 685 69006
Email: michael.scharf@ikr.uni-stuttgart.de
URI: <http://www.ikr.uni-stuttgart.de/en/~scharf>

Sally Floyd
ICIR (ICSI Center for Internet Research)

Phone: +1 (510) 666-2989
Email: floyd@icir.org
URI: <http://www.icir.org/floyd/>

Pasi Sarolahti
Nokia Research Center
P.O. Box 407
FI-00045 NOKIA GROUP
Finland

Phone: +358 50 4876607
Email: pasi.sarolahti@iki.fi

Full Copyright Statement

Copyright (C) The IETF Trust (2007).

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.

Acknowledgment

Funding for the RFC Editor function is provided by the IETF Administrative Support Activity (IASA).