

Network Working Group
INTERNET-DRAFT
Intended Status: Informational
Expires: November 16, 2013

W. Zhu
Google, Inc.
M. Jennings
May 15, 2013

Implications of Full-Duplex HTTP
draft-zhu-http-fullduplex-07.txt

Status of this Memo

This Internet-Draft is submitted to IETF in full conformance with the provisions of [BCP 78](#) and [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/lid-abstracts.html>

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

Copyright and License Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to [BCP 78](#) and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Abstract

Full-duplex HTTP follows the basic HTTP request-response semantics but also allows the server to send response body to the client at the same time when the client is transmitting request body to the server. Requirements for full-duplex HTTP are under-specified in the existing HTTP specification [[RFC2616](#)], and this memo intends to clarify the requirements of full-duplex HTTP on top of the basic HTTP protocol semantics.

Table of Contents

1	Introduction	3
1.1	Terminology	3
2	Streaming in HTTP	3
2.1	Request and Response Streaming	3
2.2	Full-Duplex Streaming	4
3	Protocol Considerations	4
3.1	Initialization	4
3.2	Termination	4
3.3	Persistent Connections	5
3.4	Time-Out	5
3.5	Proxies	5
3.6	Errors	5
4	Application Considerations	6
4.1	Compatibility	6
4.2	Fall Back	6
4.3	Buffering	6
4.4	Framing	6
5	Security Considerations	8
6	IANA Considerations	8
7	Acknowledgments	8
8	Normative References	8
	Author's Addresses	8

1 Introduction

HTTP [[RFC2616](#)] is a stateless, RPC style protocol which requires the communication between client and server follows a strict request-response pattern.

HTTP may also be used to stream data from either client or server. When bi-directional streaming is required, two TCP connections are often employed to stream client and server data separately. Using two separate connections not only introduces overhead, but also makes HTTP insufficient to be used as a standalone protocol, i.e. an application-level overlay protocol has to be defined to manage the two simultaneous connections.

However, if the server is allowed to send the response to the client at the same time that a request is being transmitted from the client to the server, then effectively full-duplex streaming becomes possible under the standard HTTP protocol semantics [[RFC2616](#)].

Full-duplex streaming requires end-to-end support from both client and server. More specifically, the client has to explicitly request for such a capability.

Given the unique properties of full-duplex HTTP, special requirements exist for both client and server. These requirements need be clearly identified so that implementations will be able to follow the expected behaviors in adopting full-duplex HTTP.

Full-duplex HTTP does not concern the underlying transport or session level protocols, e.g. those developed for SPDY or HTTP/2.0.

1.1 Terminology

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](#) [[RFC2119](#)].

2 Streaming in HTTP

2.1 Request and Response Streaming

Request streaming requires the server to deliver to the application in real-time any streaming data that has been received by the server.

Response streaming requires the server to send to the client in real-time any streaming data that has been generated by the application.

Chunked transfer-encoding [[RFC2616](#)] is expected in either request or response streaming. However, client or server should not design the underlying streaming or messaging API based on the chunked transfer-encoding as transfer-encoding is applied only hop-by-hop.

Most browsers don't yet support chunked requests, although the latest W3C XHR2 spec added streaming APIs.

[2.2](#) Full-Duplex Streaming

When full-duplex streaming is enabled, request and response are transmitted between client and server simultaneously over the same TCP connection.

Full-duplex streaming may be applied to any response that is designed to concurrently stream request and response data, for example a voice translator.

Full-duplex streaming should still follow the standard HTTP request-response semantics and maintain the basic temporal dependency between the request and response:

- 1 A server SHOULD NOT generate any response until it has received a new request. That is, unsolicited server-initiated response is not allowed.
- 2 A server SHOULD NOT complete the response prematurely when the client has not completed the request. That is, a server can only complete the response after it has received the entire request body from the client.

[3](#) Protocol Considerations

[3.1](#) Initialization

Full-duplex streaming may be started as soon as the server receives the first byte of the request body. This behavior is significantly different from the technique commonly-known as Hanging GET, which generates a separate GET request for every single piece of data that is transmitted to the client.

[3.2](#) Termination

A client terminates request streaming by completing the request, i.e. sending out the last-chunk [[RFC2616](#)] of request. The server may choose to continue streaming the response after it receives the complete request from the client, as decided by the application.

Eventually the server terminates response streaming by sending out the last-chunk of response.

Before a server receives the complete request from the client, the only way for the server to terminate response streaming is to close the connection. It is considered an illegal state for an HTTP connection to have a pending request when the response has already been completed. A client should close the connection immediately if it receives the last-chunk of response when it is still streaming the request.

[3.3](#) Persistent Connections

It's important that HTTP keep-alive and pipelining still work with full-duplex HTTP. The client and server should respect the standard semantics of HTTP persistent connections [[RFC2616](#)].

Request streaming will make the connection unavailable for pipelined (GET) requests. Also, continued (POST) response streaming after the request has been completed will prevent any pipelined requests from being transmitted over the same TCP connection.

[3.4](#) Time-Out

HTTP server or client may time-out connections while waiting for request or response. Full-duplex HTTP should not override this behavior. A connection may be closed due to time-out when either request or response streaming becomes inactive for too long.

Full-duplex HTTP introduces no special requirements on time-out of the underlying TCP connection. When time-out does happen both request and response streaming will be terminated.

[3.5](#) Proxies

HTTP proxies may not support concurrent responses, and one of the purposes of this document is to increase the awareness of full-duplex HTTP communication in proxies.

Proxies may also buffer streamed requests or responses, or have problems to handle chunked transfer-encoding for requests.

[3.6](#) Errors

According to [RFC2616](#) [[RFC2616](#)], a client should close the connection if it receives an error response when it is still transmitting the request. Full-duplex HTTP must respect this requirement.

When a server is incapable of streaming response or decides to time-out, it should lose the connection. This is also true for the client when it is streaming request data.

A client or server should stop streaming any new data after it notices that the underlying TCP connection has been closed by the other party. In-flight data will be discarded under such a half-close behavior.

4 Application Considerations

4.1 Compatibility

Full-duplex streaming is completely controlled by the server application, and should only be enabled for those clients that have been explicitly identified by the server.

It is not sufficient to enable full-duplex HTTP solely based on the User-Agent information.

For non-controlled client applications, the client may advise its capability of full-duplex streaming via URL parameters or headers (for example, "X-Accept-Streaming: full-duplex;timeout=30"). Otherwise, full-duplex streaming should be disabled, or the server should return a 404 (Not Found) status if full-duplex streaming is mandatory for the requested resource.

4.2 Fall Back

Proxies may disallow early responses, or buffer requests or responses. In such a case, the application may have to switch to an alternative protocol that either uses two TCP connections or relies on various polling techniques.

One efficient way to trigger a fall back will be for the client to wait for initial response data for a short time-out period. Generally there is no reliable way for a server to distinguish between ill-behaved clients and non-compatible proxies.

4.3 Buffering

Full-duplex HTTP expects minimized buffering on either client or server. However, applications may choose to buffer a certain amount of streaming data for optimization or application-specific purposes.

4.4 Framing

Since chunked transfer-encoding isn't a reliable way to implement message-level framing, such support needs be implemented at the application layer. Alternative protocols or client-side APIs such as WebSocket may be considered if a messaging layer is required although in most cases application-level protocols are still necessary for either fallback reasons or end-to-end delivery guarantees.

5 Security Considerations

Full-duplex HTTP introduces no new security concerns beyond those known with regular HTTP communication.

6 IANA Considerations

This document does not require any actions by the IANA.

7 Acknowledgments

Thanks to Henrik Nordstrom, Jamie Lokier, and Mark Nottingham for their feedback when the subject of this document was originally brought up on [httpbis](http://httpbis.org).

8 Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.
- [RFC2616] Fielding, R., Gettys, J., Mogul, J., Frystyk, H., Masinter, L., Leach, P., and T. Berners-Lee, "Hypertext Transfer Protocol -- HTTP/1.1", [RFC 2616](#), June 1999.

Author's Addresses

Wenbo Zhu
Google, Inc.
1600 Amphitheatre Parkway
Mountain View, CA 94043
US

Phone: +1 650 214 5894
Email: wenboz@google.com

Mike Jennings

Email: mike.c.jennings@gmail.com

