TAPS Internet-Draft Intended status: Informational Expires: September 1, 2018 M. Welzl S. Gjessing University of Oslo February 28, 2018

A Minimal Set of Transport Services for TAPS Systems draft-ietf-taps-minset-02

Abstract

This draft recommends a minimal set of IETF Transport Services offered by end systems supporting TAPS, and gives guidance on choosing among the available mechanisms and protocols. It is based on the set of transport features in <u>RFC 8303</u>.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of <u>BCP 78</u> and <u>BCP 79</u>.

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <u>https://datatracker.ietf.org/drafts/current/</u>.

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."

This Internet-Draft will expire on September 1, 2018.

Copyright Notice

Copyright (c) 2018 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to <u>BCP 78</u> and the IETF Trust's Legal Provisions Relating to IETF Documents (<u>https://trustee.ietf.org/license-info</u>) 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.

Table of Contents

	2
<u>2</u> . Terminology	<u>4</u>
$\underline{3}$. The Minimal Set of Transport Features	<u>5</u>
3.1. ESTABLISHMENT, AVAILABILITY and TERMINATION	<u>5</u>
3.2. MAINTENANCE	<u>8</u>
<u>3.2.1</u> . Connection groups	<u>8</u>
<u>3.2.2</u> . Individual connections	<u>10</u>
<u>3.3</u> . DATA Transfer	<u>10</u>
<u>3.3.1</u> . Sending Data	<u>10</u>
<u>3.3.2</u> . Receiving Data	<u>11</u>
$\underline{4}$. Conclusion	<u>12</u>
5. Acknowledgements	<u>12</u>
<u>6</u> . IANA Considerations	<u>12</u>
<u>7</u> . Security Considerations	<u>12</u>
<u>8</u> . References	<u>13</u>
<u>8.1</u> . Normative References	<u>13</u>
<u>8.2</u> . Informative References	<u>13</u>
Appendix A. Deriving the minimal set	<u>15</u>
A.1. Step 1: Categorization The Superset of Transport	
Features	<u>15</u>
Features	<u>15</u> <u>17</u>
FeaturesFeature	<u>15</u> <u>17</u> <u>32</u>
FeaturesFeature	<u>15</u> <u>17</u> <u>32</u>
Features	15 17 32 37
Features	<u>15</u> <u>17</u> <u>32</u> <u>37</u> <u>38</u>
Features	15 17 32 37 38 39
Features	15 17 32 37 38 39 40
Features	15 17 32 37 38 39 40 40
Features	15 17 32 37 38 39 40 40 41
Features	15 17 32 37 38 39 40 40 41 42
Features	15 17 32 37 38 39 40 40 41 42 43
Features	15 17 32 37 38 39 40 41 42 43
Features	15 17 32 37 38 39 40 40 41 42 43 43 43
Features	15 17 32 37 38 39 40 40 40 41 42 43 43 44 44
Features	15 17 32 37 38 39 40 40 40 41 42 43 43 44 44 45

1. Introduction

The task of any system that implements TAPS is to offer transport services to its applications, i.e. the applications running on top of the transport system, without binding them to a particular transport protocol. Currently, the set of transport services that most applications use is based on TCP and UDP (and protocols that are layered on top of them); this limits the ability for the network stack to make use of features of other transport protocols. For example, if a protocol supports out-of-order message delivery but applications always assume that the network provides an ordered bytestream, then the network stack can not immediately deliver a message that arrives out-of-order: doing so would break a fundamental assumption of the application. The net result is unnecessary headof-line blocking delay.

By exposing the transport services of multiple transport protocols, a TAPS transport system can make it possible to use these services without having to statically bind an application to a specific transport protocol. The first step towards the design of such a system was taken by [<u>RFC8095</u>], which surveys a large number of transports, and [<u>RFC8303</u>] as well as [<u>RFC8304</u>], which identify the specific transport features that are exposed to applications by the protocols TCP, MPTCP, UDP(-Lite) and SCTP as well as the LEDBAT congestion control mechanism. This memo is based on these documents and follows the same terminology (also listed below). Because the considered transport protocols conjointly cover a wide range of transport features, there is reason to hope that the resulting set (and the reasoning that led to it) will also apply to many aspects of other transport protocols.

The number of transport features of current IETF transports is large, and exposing all of them has a number of disadvantages: generally, the more functionality is exposed, the less freedom a transport system has to automate usage of the various functions of its available set of transport protocols. Some functions only exist in one particular protocol, and if an application would use them, this would statically tie the application to this protocol, counteracting the purpose of TAPS. Also, if the number of exposed features is exceedingly large, a transport system might become very difficult to use for an application programmer. Taking [RFC8303] as a basis, this document therefore develops a minimal set of transport features, removing the ones that could be harmful to the purpose of TAPS but keeping the ones that must be retained for applications to benefit from useful transport functionality.

Applications use a wide variety of APIs today. The transport features in the minimal set in this document must be reflected in *all* network APIs in order for the underlying functionality to become usable everywhere. For example, it does not help an application that talks to a middleware if only the Berkeley Sockets API is extended to offer "unordered message delivery", but the middleware only offers an ordered bytestream. Both the Berkeley Sockets API and the middleware would have to expose the "unordered message delivery" transport feature (alternatively, there may be ways for certain types of middleware to use this transport feature without exposing it, based on knowledge about the applications -- but this is

not the general case). In most situations, in the interest of being as flexible and efficient as possible, the best choice will be for a middleware or library to expose at least all of the transport features that are recommended as a "minimal set" here.

This "minimal set" can be implemented one-sided over TCP (or UDP, if certain limitations are put in place). This means that a sender-side TAPS system implementing it can talk to a non-TAPS TCP (or UDP) receiver, and a receiver-side TAPS system implementing it can talk to a non-TAPS TCP (or UDP) sender.

2. Terminology

The following terms are used throughout this document, and in subsequent documents produced by TAPS that describe the composition and decomposition of transport services.

- Transport Feature: a specific end-to-end feature that the transport layer provides to an application. Examples include confidentiality, reliable delivery, ordered delivery, messageversus-stream orientation, etc.
- Transport Service: a set of Transport Features, without an association to any given framing protocol, which provides a complete service to an application.
- Transport Protocol: an implementation that provides one or more different transport services using a specific framing and header format on the wire.
- Transport Service Instance: an arrangement of transport protocols with a selected set of features and configuration parameters that implements a single transport service, e.g., a protocol stack (RTP over UDP).
- Application: an entity that uses the transport layer for end-to-end delivery data across the network (this may also be an upper layer protocol or tunnel encapsulation).
- Application-specific knowledge: knowledge that only applications have.
- Endpoint: an entity that communicates with one or more other endpoints using a transport protocol.
- Connection: shared state of two or more endpoints that persists across messages that are transmitted between these endpoints.
- Socket: the combination of a destination IP address and a destination port number.

Moreover, throughout the document, the protocol name "UDP(-Lite)" is used when discussing transport features that are equivalent for UDP and UDP-Lite; similarly, the protocol name "TCP" refers to both TCP and MPTCP.

3. The Minimal Set of Transport Features

Based on the categorization, reduction and discussion in <u>Appendix A</u>, this section describes the minimal set of transport features that is offered by end systems supporting TAPS. The described transport system can be implemented over TCP; elements of the system that may prohibit implementation over UDP are marked with "!UDP". To implement a transport system that can also work over UDP, these marked transport features should be excluded.

As in <u>Appendix A</u>, <u>Appendix A.2</u> and [<u>RFC8303</u>], we categorize the minimal set of transport features as 1) CONNECTION related (ESTABLISHMENT, AVAILABILITY, MAINTENANCE, TERMINATION) and 2) DATA Transfer related (Sending Data, Receiving Data, Errors). Here, the focus is on connections that the transport system offers, as opposed to connections of transport protocols that the transport system uses.

<u>3.1</u>. ESTABLISHMENT, AVAILABILITY and TERMINATION

A connection must first be "created" to allow for some initial configuration to be carried out before the transport system can actively or passively establish communication with a remote endpoint. All configuration parameters in <u>Section 3.2</u> can be used initially, although some of them may only take effect when a connection has been established with a chosen transport protocol. Configuring a connection early helps a transport system make the right decisions. For example, grouping information can influence the transport system to implement a connection as a stream of a multi-streaming protocol's existing association or not.

For ungrouped connections, early configuration is necessary because it allows the transport system to know which protocols it should try to use (to steer a mechanism such as "Happy Eyeballs" [I-D.grinnemo-taps-he]). In particular, a transport system that only makes a one-time choice for a particular protocol must know early about strict requirements that must be kept, or it can end up in a deadlock situation (e.g., having chosen UDP and later be asked to support reliable transfer). As a possibility to correctly handle these cases, we provide the following decision tree (this is derived from <u>Appendix A.2.1</u> excluding authentication, as explained in <u>Section 7</u>):

```
Minimal TAPS Transport Services February 2018
- Will it ever be necessary to offer any of the following?
  * Reliably transfer data
  * Notify the peer of closing/aborting
  * Preserve data ordering
 Yes: SCTP or TCP can be used.
  - Is any of the following useful to the application?
    * Choosing a scheduler to operate between connections
     in a group, with the possibility to configure a priority
     or weight per connection
    * Configurable message reliability
   * Unordered message delivery
   * Request not to delay the acknowledgement (SACK) of a message
   Yes: SCTP is preferred.
   No:
    - Is any of the following useful to the application?
      * Hand over a message to reliably transfer (possibly
       multiple times) before connection establishment
      * Suggest timeout to the peer
      * Notification of Excessive Retransmissions (early
       warning below abortion threshold)
      * Notification of ICMP error message arrival
     Yes: TCP is preferred.
     No: SCTP and TCP are equally preferable.
 No: all protocols can be used.
  - Is any of the following useful to the application?
   * Specify checksum coverage used by the sender
   * Specify minimum checksum coverage required by receiver
   Yes: UDP-Lite is preferred.
   No: UDP is preferred.
Note that this decision tree is not optimal for all cases. For
example, if an application wants to use "Specify checksum coverage
```

used by the sender", which is only offered by UDP-Lite, and "Configure priority or weight for a scheduler", which is only offered by SCTP, the above decision tree will always choose UDP-Lite, making it impossible to use SCTP's schedulers with priorities between grouped connections. The transport system must know which choice is more important for the application in order to make the best decision. We caution implementers to be aware of the full set of trade-offs, for which we recommend consulting the list in Appendix A.2.1 when deciding how to initialize a connection.

To summarize, the following parameters serve as input for the transport system to help it choose and configure a suitable protocol:

- o Reliability: a boolean that should be set to true when any of the following will be useful to the application: reliably transfer data; notify the peer of closing/aborting; preserve data ordering.
- Checksum_coverage: a boolean to specify whether it will be useful to the application to specify checksum coverage when sending or receiving.
- o Config_msg_prio: a boolean that should be set to true when any of the following per-message configuration or prioritization mechanisms will be useful to the application: choosing a scheduler to operate between grouped connections, with the possibility to configure a priority or weight per connection; configurable message reliability; unordered message delivery; requesting not to delay the acknowledgement (SACK) of a message.
- o Earlymsg_timeout_notifications: a boolean that should be set to true when any of the following will be useful to the application: hand over a message to reliably transfer (possibly multiple times) before connection establishment; suggest timeout to the peer; notification of excessive retransmissions (early warning below abortion threshold); notification of ICMP error message arrival.

Once a connection is created, it can be queried for the maximum amount of data that an application can possibly expect to have reliably transmitted before or during transport connection establishment (with zero being a possible answer) (see <u>Section 3.2.1</u>). An application can also give the connection a message for reliable transmission before or during connection establishment (!UDP); the transport system will then try to transmit it as early as possible. An application can facilitate sending a message particularly early by marking it as "idempotent" (see <u>Section 3.3.1</u>); in this case, the receiving application must be prepared to potentially receive multiple copies of the message (because idempotent messages are reliably transferred, asking for idempotence is not necessary for systems that support UDP).

After creation, a transport system can actively establish communication with a peer, or it can passively listen for incoming connection requests. Note that active establishment may or may not trigger a notification on the listening side. It is possible that the first notification on the listening side is the arrival of the first data that the active side sends (a receiver-side transport system could handle this by continuing to block a "Listen" call, immediately followed by issuing "Receive", for example; callbackbased implementations could simply skip the equivalent of "Listen"). This also means that the active opening side is assumed to be the first side sending data. A transport system can actively close a connection, i.e. terminate it after reliably delivering all remaining data to the peer (if reliable data delivery was requested earlier (!UDP)), in which case the peer is notified that the connection is closed. Alternatively, a connection can be aborted without delivering outstanding data to the peer. In case reliable or partially reliable data delivery was requested earlier (!UDP), the peer is notified that the connection is aborted. A timeout can be configured to abort a connection when data could not be delivered for too long (!UDP); however, timeout-based abortion does not notify the peer application that the connection has been aborted. Because half-closed connections are not supported, when a host implementing TAPS receives a notification that the peer is closing or aborting the connection (!UDP), its peer may not be able to read outstanding data. This means that unacknowledged data residing a transport system's send buffer may have to be dropped from that buffer upon arrival of a "close" or "abort" notification from the peer.

3.2. MAINTENANCE

A transport system must offer means to group connections, but it cannot guarantee truly grouping them using the transport protocols that it uses (e.g., it cannot be guaranteed that connections become multiplexed as streams on a single SCTP association when SCTP may not be available). The transport system must therefore ensure that group- versus non-group-configurations are handled correctly in some way (e.g., by applying the configuration to all grouped connections even when they are not multiplexed, or informing the application about grouping success or failure).

As a general rule, any configuration described below should be carried out as early as possible to aid the transport system's decision making.

3.2.1. Connection groups

The following transport features and notifications (some directly from Appendix A.2, some new or changed, based on the discussion in Appendix A.3) automatically apply to all grouped connections:

(!UDP) Configure a timeout: this can be done with the following parameters:

- o A timeout value for aborting connections, in seconds
- o A timeout value to be suggested to the peer (if possible), in seconds
- o The number of retransmissions after which the application should be notifed of "Excessive Retransmissions"

Configure urgency: this can be done with the following parameters:

- o A number to identify the type of scheduler that should be used to operate between connections in the group (no guarantees given). Schedulers are defined in [RFC8260].
- o A "capacity profile" number to identify how an application wants to use its available capacity. Choices can be "lowest possible latency at the expense of overhead" (which would disable any Nagle-like algorithm), "scavenger", or values that help determine the DSCP value for a connection (e.g. similar to table 1 in [I-D.ietf-tsvwg-rtcweb-gos]).
- o A buffer limit (in bytes); when the sender has less then the provided limit of bytes in the buffer, the application may be notified. Notifications are not guaranteed, and it is optional for a transport system to support buffer limit values greater than 0. Note that this limit and its notification should operate across the buffers of the whole transport system, i.e. also any potential buffers that the transport system itself may use on top of the transport's send buffer.

Following Appendix A.3.7, these properties can be queried:

- o The maximum message size that may be sent without fragmentation via the configured interface. This is optional for a transport system to offer, and may return an error ("not available"). It can aid applications implementing Path MTU Discovery.
- o The maximum transport message size that can be sent, in bytes. Irrespective of fragmentation, there is a size limit for the messages that can be handed over to SCTP or UDP(-Lite); because the service provided by a transport system is independent of the transport protocol, it must allow an application to query this value -- the maximum size of a message in an Application-Framed-Bytestream (see Appendix A.3.1). This may also return an error when data is not delimited ("not available").
- o The maximum transport message size that can be received from the configured interface, in bytes (or "not available").
- o The maximum amount of data that can possibly be sent before or during connection establishment, in bytes.

In addition to the already mentioned closing / aborting notifications and possible send errors, the following notifications can occur:

- o Excessive Retransmissions: the configured (or a default) number of retransmissions has been reached, yielding this early warning below an abortion threshold.
- o ICMP Arrival (parameter: ICMP message): an ICMP packet carrying the conveyed ICMP message has arrived.

- o ECN Arrival (parameter: ECN value): a packet carrying the conveyed ECN value has arrived. This can be useful for applications implementing congestion control.
- o Timeout (parameter: s seconds): data could not be delivered for s seconds.
- o Drain: the send buffer has either drained below the configured buffer limit or it has become completely empty. This is a generic notification that tries to enable uniform access to "TCP_NOTSENT_LOWAT" as well as the "SENDER DRY" notification (as discussed in Appendix A.3.4 -- SCTP's "SENDER DRY" is a special case where the threshold (for unsent data) is 0 and there is also no more unacknowledged data in the send buffer).

3.2.2. Individual connections

Configure priority or weight for a scheduler, as described in [<u>RFC8260</u>].

Configure checksum usage: this can be done with the following parameters, but there is no guarantee that any checksum limitations will indeed be enforced (the default behavior is "full coverage, checksum enabled"):

- o A boolean to enable / disable usage of a checksum when sending
- o The desired coverage (in bytes) of the checksum used when sending
- o A boolean to enable / disable requiring a checksum when receiving
- o The required minimum coverage (in bytes) of the checksum when receiving

3.3. DATA Transfer

3.3.1. Sending Data

When sending a message, no guarantees are given about the preservation of message boundaries to the peer; if message boundaries are needed, the receiving application at the peer must know about them beforehand (or the transport system cannot use TCP). Note that an application should already be able to hand over data before the transport system establishes a connection with a chosen transport protocol. Regarding the message that is being handed over, the following parameters can be used:

o Reliability: this parameter is used to convey a choice of: fully reliable (!UDP), unreliable without congestion control, unreliable (!UDP), partially reliable (see [RFC3758] and [RFC7496] for details on how to specify partial reliability) (!UDP). The latter two choices are optional for a transport system to offer and may result in full reliability. Note that applications sending

unreliable data without congestion control should themselves perform congestion control in accordance with [<u>RFC2914</u>].

- o (!UDP) Ordered: this boolean parameter lets an application choose between ordered message delivery (true) and possibly unordered, potentially faster message delivery (false).
- o Bundle: a boolean that expresses a preference for allowing to bundle messages (true) or not (false). No guarantees are given.
- o DelAck: a boolean that, if false, lets an application request that the peer would not delay the acknowledgement for this message.
- o Fragment: a boolean that expresses a preference for allowing to fragment messages (true) or not (false), at the IP level. No guarantees are given.
- o (!UDP) Idempotent: a boolean that expresses whether a message is idempotent (true) or not (false). Idempotent messages may arrive multiple times at the receiver (but they will arrive at least once). When data is idempotent it can be used by the receiver immediately on a connection establishment attempt. Thus, if data is handed over before the transport system establishes a connection with a chosen transport protocol, stating that a message is idempotent facilitates transmitting it to the peer application particularly early.

An application can be notified of a failure to send a specific message. There is no guarantee of such notifications, i.e. send failures can also silently occur.

3.3.2. Receiving Data

A receiving application obtains an "Application-Framed Bytestream" (AFra-Bytestream); this concept is further described in <u>Appendix A.3.1</u>). In line with TCP's receiver semantics, an AFra-Bytestream is just a stream of bytes to the receiver. If message boundaries were specified by the sender, a receiver-side transport system implementing only the minimum set of transport services defined here will still not inform the receiving application about them (this limitation is only needed for transport systems that are implemented to directly use TCP).

Different from TCP's semantics, if the sending application has allowed that messages are not fully reliably transferred, or delivered out of order, then such re-ordering or unreliability may be reflected per message in the arriving data. Messages will always stay intact - i.e. if an incomplete message is contained at the end of the arriving data block, this message is guaranteed to continue in the next arriving data block.

4. Conclusion

By decoupling applications from transport protocols, a TAPS transport system provides a different abstraction level than the Berkeley sockets interface. As with high- vs. low-level programming languages, a higher abstraction level allows more freedom for automation below the interface, yet it takes some control away from the application programmer. This is the design trade-off that a transport system developer is facing, and this document provides quidance on the design of this abstraction level. Some transport features are currently rarely offered by APIs, yet they must be offered or they can never be used ("functional" transport features). Other transport features are offered by the APIs of the protocols covered here, but not exposing them in a TAPS API would allow for more freedom to automate protocol usage in a transport system. The minimal set presented in this document is an effort to find a middle ground that can be recommended for transport systems to implement, on the basis of the transport features discussed in [RFC8303].

5. Acknowledgements

The authors would like to thank all the participants of the TAPS Working Group and the NEAT and MAMI research projects for valuable input to this document. We especially thank Michael Tuexen for help with connection connection establishment/teardown and Gorry Fairhurst for his suggestions regarding fragmentation and packet sizes. This work has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No. 644334 (NEAT).

6. IANA Considerations

XX RFC ED - PLEASE REMOVE THIS SECTION XXX

This memo includes no request to IANA.

7. Security Considerations

Authentication, confidentiality protection, and integrity protection are identified as transport features by [RFC8095]. As currently deployed in the Internet, these features are generally provided by a protocol or layer on top of the transport protocol; no current fullfeatured standards-track transport protocol provides all of these transport features on its own. Therefore, these transport features are not considered in this document, with the exception of native authentication capabilities of TCP and SCTP for which the security considerations in [RFC5925] and [RFC4895] apply. Security is

discussed further in a separate TAPS document [<u>I-D.pauly-taps-transport-security</u>].

8. References

8.1. Normative References

[RFC8303] Welzl, M., Tuexen, M., and N. Khademi, "On the Usage of Transport Features Provided by IETF Transport Protocols", <u>RFC 8303</u>, DOI 10.17487/RFC8303, February 2018, <https://www.rfc-editor.org/info/rfc8303>.

8.2. Informative References

[COBS] Cheshire, S. and M. Baker, "Consistent Overhead Byte Stuffing", September 1997, <<u>http://stuartcheshire.org/papers/COBSforToN.pdf</u>>.

[I-D.grinnemo-taps-he]

Grinnemo, K., Brunstrom, A., Hurtig, P., Khademi, N., and Z. Bozakov, "Happy Eyeballs for Transport Selection", <u>draft-grinnemo-taps-he-03</u> (work in progress), July 2017.

[I-D.ietf-tsvwg-rtcweb-qos]

Jones, P., Dhesikan, S., Jennings, C., and D. Druta, "DSCP Packet Markings for WebRTC QoS", <u>draft-ietf-tsvwg-rtcweb-</u> <u>qos-18</u> (work in progress), August 2016.

[I-D.pauly-taps-transport-security]

Pauly, T., Rose, K., and C. Wood, "A Survey of Transport Security Protocols", <u>draft-pauly-taps-transport-</u> <u>security-01</u> (work in progress), January 2018.

[LBE-draft]

- Bless, R., "A Lower Effort Per-Hop Behavior (LE PHB)", Internet-draft <u>draft-tsvwg-le-phb-03</u>, February 2018.
- [RFC2914] Floyd, S., "Congestion Control Principles", <u>BCP 41</u>, <u>RFC 2914</u>, DOI 10.17487/RFC2914, September 2000, <<u>https://www.rfc-editor.org/info/rfc2914</u>>.
- [RFC3758] Stewart, R., Ramalho, M., Xie, Q., Tuexen, M., and P. Conrad, "Stream Control Transmission Protocol (SCTP) Partial Reliability Extension", <u>RFC 3758</u>, DOI 10.17487/RFC3758, May 2004, <<u>https://www.rfc-editor.org/info/rfc3758</u>>.

- [RFC4895] Tuexen, M., Stewart, R., Lei, P., and E. Rescorla, "Authenticated Chunks for the Stream Control Transmission Protocol (SCTP)", <u>RFC 4895</u>, DOI 10.17487/RFC4895, August 2007, <<u>https://www.rfc-editor.org/info/rfc4895</u>>.
- [RFC4987] Eddy, W., "TCP SYN Flooding Attacks and Common Mitigations", <u>RFC 4987</u>, DOI 10.17487/RFC4987, August 2007, <<u>https://www.rfc-editor.org/info/rfc4987</u>>.
- [RFC5925] Touch, J., Mankin, A., and R. Bonica, "The TCP Authentication Option", <u>RFC 5925</u>, DOI 10.17487/RFC5925, June 2010, <<u>https://www.rfc-editor.org/info/rfc5925</u>>.
- [RFC7305] Lear, E., Ed., "Report from the IAB Workshop on Internet Technology Adoption and Transition (ITAT)", <u>RFC 7305</u>, DOI 10.17487/RFC7305, July 2014, <<u>https://www.rfc-editor.org/info/rfc7305</u>>.
- [RFC7413] Cheng, Y., Chu, J., Radhakrishnan, S., and A. Jain, "TCP Fast Open", <u>RFC 7413</u>, DOI 10.17487/RFC7413, December 2014, <<u>https://www.rfc-editor.org/info/rfc7413</u>>.
- [RFC7496] Tuexen, M., Seggelmann, R., Stewart, R., and S. Loreto, "Additional Policies for the Partially Reliable Stream Control Transmission Protocol Extension", <u>RFC 7496</u>, DOI 10.17487/RFC7496, April 2015, <<u>https://www.rfc-editor.org/info/rfc7496</u>>.
- [RFC8095] Fairhurst, G., Ed., Trammell, B., Ed., and M. Kuehlewind, Ed., "Services Provided by IETF Transport Protocols and Congestion Control Mechanisms", <u>RFC 8095</u>, DOI 10.17487/RFC8095, March 2017, <<u>https://www.rfc-editor.org/info/rfc8095</u>>.
- [RFC8260] Stewart, R., Tuexen, M., Loreto, S., and R. Seggelmann, "Stream Schedulers and User Message Interleaving for the Stream Control Transmission Protocol", <u>RFC 8260</u>, DOI 10.17487/RFC8260, November 2017, <<u>https://www.rfc-editor.org/info/rfc8260</u>>.
- [RFC8304] Fairhurst, G. and T. Jones, "Transport Features of the User Datagram Protocol (UDP) and Lightweight UDP (UDP-Lite)", <u>RFC 8304</u>, DOI 10.17487/RFC8304, February 2018, <<u>https://www.rfc-editor.org/info/rfc8304</u>>.

[WWDC2015]

Lakhera, P. and S. Cheshire, "Your App and Next Generation Networks", Apple Worldwide Developers Conference 2015, San Francisco, USA, June 2015, <https://developer.apple.com/videos/wwdc/2015/?id=719>.

Appendix A. Deriving the minimal set

We approach the construction of a minimal set of transport features in the following way:

- 1. Categorization: the superset of transport features from [RFC8303] is presented, and transport features are categorized for later reduction.
- 2. Reduction: a shorter list of transport features is derived from the categorization in the first step. This removes all transport features that do not require application-specific knowledge or cannot be implemented with TCP or UDP.
- 3. Discussion: the resulting list shows a number of peculiarities that are discussed, to provide a basis for constructing the minimal set.
- 4. Construction: Based on the reduced set and the discussion of the transport features therein, a minimal set is constructed.

The first three steps as well as the underlying rationale for constructing the minimal set are described in this appendix. The minimal set itself is described in Section 3.

A.1. Step 1: Categorization -- The Superset of Transport Features

Following [RFC8303], we divide the transport features into two main groups as follows:

- 1. CONNECTION related transport features
 - ESTABLISHMENT
 - AVAILABILITY
 - MAINTENANCE
 - TERMINATION
- 2. DATA Transfer related transport features
 - Sending Data
 - Receiving Data
 - Errors

We assume that applications have no specific requirements that need knowledge about the network, e.g. regarding the choice of network interface or the end-to-end path. Even with these assumptions, there

Minimal TAPS Transport Services February 2018

are certain requirements that are strictly kept by transport protocols today, and these must also be kept by a transport system. Some of these requirements relate to transport features that we call "Functional".

Functional transport features provide functionality that cannot be used without the application knowing about them, or else they violate assumptions that might cause the application to fail. For example, ordered message delivery is a functional transport feature: it cannot be configured without the application knowing about it because the application's assumption could be that messages always arrive in order. Failure includes any change of the application behavior that is not performance oriented, e.g. security.

"Change DSCP" and "Disable Nagle algorithm" are examples of transport features that we call "Optimizing": if a transport system autonomously decides to enable or disable them, an application will not fail, but a transport system may be able to communicate more efficiently if the application is in control of this optimizing transport feature. These transport features require applicationspecific knowledge (e.g., about delay/bandwidth requirements or the length of future data blocks that are to be transmitted).

The transport features of IETF transport protocols that do not require application-specific knowledge and could therefore be transparently utilized by a transport system are called "Automatable".

Finally, some transport features are aggregated and/or slightly changed in the description below. These transport features are marked as "ADDED". The corresponding transport features are automatable, and they are listed immediately below the "ADDED" transport feature.

In this description, transport services are presented following the nomenclature "CATEGORY.[SUBCATEGORY].SERVICENAME.PROTOCOL", equivalent to "pass 2" in [RFC8303]. We also sketch how some of the TAPS transport features can be implemented by a transport system. For all transport features that are categorized as "functional" or "optimizing", and for which no matching TCP and/or UDP primitive exists in "pass 2" of [<u>RFC8303</u>], a brief discussion on how to implement them over TCP and/or UDP is included.

We designate some transport features as "automatable" on the basis of a broader decision that affects multiple transport features:

o Most transport features that are related to multi-streaming were designated as "automatable". This was done because the decision

on whether to use multi-streaming or not does not depend on application-specific knowledge. This means that a connection that is exhibited to an application could be implemented by using a single stream of an SCTP association instead of mapping it to a complete SCTP association or TCP connection. This could be achieved by using more than one stream when an SCTP association is first established (CONNECT.SCTP parameter "outbound stream count"), maintaining an internal stream number, and using this stream number when sending data (SEND.SCTP parameter "stream number"). Closing or aborting a connection could then simply free the stream number for future use. This is discussed further in Appendix A.3.2.

o All transport features that are related to using multiple paths or the choice of the network interface were designated as "automatable". Choosing a path or an interface does not depend on application-specific knowledge. For example, "Listen" could always listen on all available interfaces and "Connect" could use the default interface for the destination IP address.

A.1.1. CONNECTION Related Transport Features

ESTABLISHMENT:

```
o Connect
  Protocols: TCP, SCTP, UDP(-Lite)
  Functional because the notion of a connection is often reflected
  in applications as an expectation to be able to communicate after
  a "Connect" succeeded, with a communication sequence relating to
  this transport feature that is defined by the application
  protocol.
```

Implementation: via CONNECT.TCP, CONNECT.SCTP or CONNECT.UDP(-Lite).

- o Specify which IP Options must always be used Protocols: TCP, UDP(-Lite) Automatable because IP Options relate to knowledge about the network, not the application.
- o Request multiple streams Protocols: SCTP Automatable because using multi-streaming does not require application-specific knowledge. Implementation: see <u>Appendix A.3.2</u>.

Internet-Draft Minimal TAPS Transport Services February 2018

o Limit the number of inbound streams Protocols: SCTP Automatable because using multi-streaming does not require application-specific knowledge. Implementation: see Appendix A.3.2.

o Specify number of attempts and/or timeout for the first establishment message Protocols: TCP, SCTP Functional because this is closely related to potentially assumed reliable data delivery for data that is sent before or during connection establishment. Implementation: Using a parameter of CONNECT.TCP and CONNECT.SCTP. Implementation over UDP: Do nothing (this is irrelevant in case of UDP because there, reliable data delivery is not assumed).

o Obtain multiple sockets Protocols: SCTP Automatable because the usage of multiple paths to communicate to the same end host relates to knowledge about the network, not the application.

o Disable MPTCP

Protocols: MPTCP

Automatable because the usage of multiple paths to communicate to the same end host relates to knowledge about the network, not the application.

Implementation: via a boolean parameter in CONNECT.MPTCP.

o Configure authentication

Protocols: TCP, SCTP

Functional because this has a direct influence on security. Implementation: via parameters in CONNECT.TCP and CONNECT.SCTP. Implementation over TCP: With TCP, this allows to configure Master Key Tuples (MKTs) to authenticate complete segments (including the TCP IPv4 pseudoheader, TCP header, and TCP data). With SCTP, this allows to specify which chunk types must always be authenticated. Authenticating only certain chunk types creates a reduced level of security that is not supported by TCP; to be compatible, this should therefore only allow to authenticate all chunk types. Key material must be provided in a way that is compatible with both [<u>RFC4895</u>] and [<u>RFC5925</u>].

Implementation over UDP: Not possible.

- o Indicate (and/or obtain upon completion) an Adaptation Layer via an adaptation code point Protocols: SCTP Functional because it allows to send extra data for the sake of identifying an adaptation layer, which by itself is applicationspecific. Implementation: via a parameter in CONNECT.SCTP. Implementation over TCP: not possible. Implementation over UDP: not possible.
- o Request to negotiate interleaving of user messages Protocols: SCTP Automatable because it requires using multiple streams, but requesting multiple streams in the CONNECTION.ESTABLISHMENT category is automatable. Implementation: via a parameter in CONNECT.SCTP.
- o Hand over a message to reliably transfer (possibly multiple times) before connection establishment Protocols: TCP Functional because this is closely tied to properties of the data that an application sends or expects to receive. Implementation: via a parameter in CONNECT.TCP. Implementation over UDP: not possible.
- o Hand over a message to reliably transfer during connection establishment Protocols: SCTP Functional because this can only work if the message is limited in size, making it closely tied to properties of the data that an application sends or expects to receive. Implementation: via a parameter in CONNECT.SCTP. Implementation over UDP: not possible.
- o Enable UDP encapsulation with a specified remote UDP port number Protocols: SCTP Automatable because UDP encapsulation relates to knowledge about the network, not the application.
Internet-Draft

AVAILABILITY:

o Listen Protocols: TCP, SCTP, UDP(-Lite) Functional because the notion of accepting connection requests is often reflected in applications as an expectation to be able to communicate after a "Listen" succeeded, with a communication sequence relating to this transport feature that is defined by the application protocol. ADDED. This differs from the 3 automatable transport features below in that it leaves the choice of interfaces for listening open. Implementation: by listening on all interfaces via LISTEN.TCP (not providing a local IP address) or LISTEN.SCTP (providing SCTP port number / address pairs for all local IP addresses). LISTEN.UDP(-Lite) supports both methods. o Listen, 1 specified local interface Protocols: TCP, SCTP, UDP(-Lite) Automatable because decisions about local interfaces relate to knowledge about the network and the Operating System, not the application. o Listen, N specified local interfaces Protocols: SCTP Automatable because decisions about local interfaces relate to knowledge about the network and the Operating System, not the application. o Listen, all local interfaces Protocols: TCP, SCTP, UDP(-Lite) Automatable because decisions about local interfaces relate to knowledge about the network and the Operating System, not the application.

o Specify which IP Options must always be used Protocols: TCP, UDP(-Lite) Automatable because IP Options relate to knowledge about the network, not the application.

- o Disable MPTCP Protocols: MPTCP Automatable because the usage of multiple paths to communicate to the same end host relates to knowledge about the network, not the application.
- o Configure authentication Protocols: TCP, SCTP Functional because this has a direct influence on security. Implementation: via parameters in LISTEN.TCP and LISTEN.SCTP. Implementation over TCP: With TCP, this allows to configure Master Key Tuples (MKTs) to authenticate complete segments (including the TCP IPv4 pseudoheader, TCP header, and TCP data). With SCTP, this allows to specify which chunk types must always be authenticated. Authenticating only certain chunk types creates a reduced level of security that is not supported by TCP; to be compatible, this should therefore only allow to authenticate all chunk types. Key material must be provided in a way that is compatible with both [RFC4895] and [RFC5925].

Implementation over UDP: not possible.

- o Obtain requested number of streams Protocols: SCTP Automatable because using multi-streaming does not require application-specific knowledge. Implementation: see <u>Appendix A.3.2</u>.
- o Limit the number of inbound streams Protocols: SCTP Automatable because using multi-streaming does not require application-specific knowledge. Implementation: see Appendix A.3.2.
- o Indicate (and/or obtain upon completion) an Adaptation Layer via an adaptation code point Protocols: SCTP Functional because it allows to send extra data for the sake of identifying an adaptation layer, which by itself is applicationspecific. Implementation: via a parameter in LISTEN.SCTP. Implementation over TCP: not possible. Implementation over UDP: not possible.

o Request to negotiate interleaving of user messages Protocols: SCTP Automatable because it requires using multiple streams, but requesting multiple streams in the CONNECTION.ESTABLISHMENT category is automatable. Implementation: via a parameter in LISTEN.SCTP.

MAINTENANCE:

o Change timeout for aborting connection (using retransmit limit or time value) Protocols: TCP, SCTP Functional because this is closely related to potentially assumed reliable data delivery. Implementation: via CHANGE_TIMEOUT.TCP or CHANGE_TIMEOUT.SCTP. Implementation over UDP: not possible (UDP is unreliable and there is no connection timeout).

- o Suggest timeout to the peer Protocols: TCP Functional because this is closely related to potentially assumed reliable data delivery. Implementation: via CHANGE_TIMEOUT.TCP. Implementation over UDP: not possible (UDP is unreliable and there is no connection timeout).
- o Disable Nagle algorithm Protocols: TCP, SCTP Optimizing because this decision depends on knowledge about the size of future data blocks and the delay between them. Implementation: via DISABLE_NAGLE.TCP and DISABLE_NAGLE.SCTP. Implementation over UDP: do nothing (UDP does not implement the Nagle algorithm).
- o Request an immediate heartbeat, returning success/failure Protocols: SCTP Automatable because this informs about network-specific knowledge.

- o Notification of Excessive Retransmissions (early warning below abortion threshold) Protocols: TCP Optimizing because it is an early warning to the application, informing it of an impending functional event. Implementation: via ERROR.TCP. Implementation over UDP: do nothing (there is no abortion threshold).
- o Add path Protocols: MPTCP, SCTP MPTCP Parameters: source-IP; source-Port; destination-IP; destination-Port SCTP Parameters: local IP address Automatable because the usage of multiple paths to communicate to the same end host relates to knowledge about the network, not the application.
- o Remove path Protocols: MPTCP, SCTP MPTCP Parameters: source-IP; source-Port; destination-IP; destination-Port SCTP Parameters: local IP address Automatable because the usage of multiple paths to communicate to the same end host relates to knowledge about the network, not the application.
- o Set primary path Protocols: SCTP Automatable because the usage of multiple paths to communicate to the same end host relates to knowledge about the network, not the application.
- o Suggest primary path to the peer Protocols: SCTP Automatable because the usage of multiple paths to communicate to the same end host relates to knowledge about the network, not the application.

- o Configure Path Switchover Protocols: SCTP Automatable because the usage of multiple paths to communicate to the same end host relates to knowledge about the network, not the application.
- o Obtain status (query or notification) Protocols: SCTP, MPTCP SCTP parameters: association connection state; destination transport address list; destination transport address reachability states; current local and peer receiver window size; current local congestion window sizes; number of unacknowledged DATA chunks; number of DATA chunks pending receipt; primary path; most recent SRTT on primary path; RTO on primary path; SRTT and RTO on other destination addresses; MTU per path; interleaving supported yes/no MPTCP parameters: subflow-list (identified by source-IP; source-Port; destination-IP; destination-Port) Automatable because these parameters relate to knowledge about the network, not the application.
- o Specify DSCP field Protocols: TCP, SCTP, UDP(-Lite) Optimizing because choosing a suitable DSCP value requires application-specific knowledge. Implementation: via SET_DSCP.TCP / SET_DSCP.SCTP / SET_DSCP.UDP(-Lite)
- o Notification of ICMP error message arrival Protocols: TCP, UDP(-Lite) Optimizing because these messages can inform about success or failure of functional transport features (e.g., host unreachable relates to "Connect") Implementation: via ERROR.TCP or ERROR.UDP(-Lite).
- o Obtain information about interleaving support Protocols: SCTP Automatable because it requires using multiple streams, but requesting multiple streams in the CONNECTION.ESTABLISHMENT category is automatable. Implementation: via STATUS.SCTP.

- o Change authentication parameters Protocols: TCP, SCTP Functional because this has a direct influence on security. Implementation: via SET_AUTH.TCP and SET_AUTH.SCTP. Implementation over TCP: With SCTP, this allows to adjust key_id, key, and hmac_id. With TCP, this allows to change the preferred outgoing MKT (current_key) and the preferred incoming MKT (rnext_key), respectively, for a segment that is sent on the connection. Key material must be provided in a way that is compatible with both [RFC4895] and [RFC5925]. Implementation over UDP: not possible.
- o Obtain authentication information Protocols: SCTP Functional because authentication decisions may have been made by the peer, and this has an influence on the necessary applicationlevel measures to provide a certain level of security. Implementation: via GET_AUTH.SCTP. Implementation over TCP: With SCTP, this allows to obtain key_id and a chunk list. With TCP, this allows to obtain current_key and rnext_key from a previously received segment. Key material must be provided in a way that is compatible with both [RFC4895] and [RFC5925].

Implementation over UDP: not possible.

- o Reset Stream Protocols: SCTP Automatable because using multi-streaming does not require application-specific knowledge. Implementation: see Appendix A.3.2.
- o Notification of Stream Reset Protocols: STCP Automatable because using multi-streaming does not require application-specific knowledge. Implementation: see <u>Appendix A.3.2</u>.
- o Reset Association Protocols: SCTP Automatable because deciding to reset an association does not require application-specific knowledge. Implementation: via RESET_ASSOC.SCTP.

Internet-Draft Min

- Notification of Association Reset
 Protocols: STCP
 Automatable because this notification does not relate to application-specific knowledge.
- Add Streams
 Protocols: SCTP
 Automatable because using multi-streaming does not require application-specific knowledge.
 Implementation: see <u>Appendix A.3.2</u>.
- Notification of Added Stream
 Protocols: STCP
 Automatable because using multi-streaming does not require application-specific knowledge.
 Implementation: see <u>Appendix A.3.2</u>.
- Choose a scheduler to operate between streams of an association Protocols: SCTP
 Optimizing because the scheduling decision requires applicationspecific knowledge. However, if a transport system would not use this, or wrongly configure it on its own, this would only affect the performance of data transfers; the outcome would still be correct within the "best effort" service model.
 Implementation: using SET_STREAM_SCHEDULER.SCTP.
 Implementation over TCP: do nothing.
 Implementation over UDP: do nothing.
- Configure priority or weight for a scheduler Protocols: SCTP
 Optimizing because the priority or weight requires applicationspecific knowledge. However, if a transport system would not use this, or wrongly configure it on its own, this would only affect the performance of data transfers; the outcome would still be correct within the "best effort" service model.
 Implementation: using CONFIGURE_STREAM_SCHEDULER.SCTP.
 Implementation over TCP: do nothing.
 Implementation over UDP: do nothing.

o Configure send buffer size

Protocols: SCTP Automatable because this decision relates to knowledge about the network and the Operating System, not the application (see also the discussion in Appendix A.3.4).

- o Configure receive buffer (and rwnd) size Protocols: SCTP Automatable because this decision relates to knowledge about the network and the Operating System, not the application.
- o Configure message fragmentation Protocols: SCTP Automatable because fragmentation relates to knowledge about the network and the Operating System, not the application. Implementation: by always enabling it with CONFIG_FRAGMENTATION.SCTP and auto-setting the fragmentation size based on network or Operating System conditions.
- o Configure PMTUD Protocols: SCTP Automatable because Path MTU Discovery relates to knowledge about the network, not the application.
- o Configure delayed SACK timer Protocols: SCTP Automatable because the receiver-side decision to delay sending SACKs relates to knowledge about the network, not the application (it can be relevant for a sending application to request not to delay the SACK of a message, but this is a different transport feature).
- o Set Cookie life value Protocols: SCTP Functional because it relates to security (possibly weakened by keeping a cookie very long) versus the time between connection establishment attempts. Knowledge about both issues can be application-specific.

Implementation over TCP: the closest specified TCP functionality is the cookie in TCP Fast Open; for this, [<u>RFC7413</u>] states that the server "can expire the cookie at any time to enhance security" and section 4.1.2 describes an example implementation where updating the key on the server side causes the cookie to expire. Alternatively, for implementations that do not support TCP Fast Open, this transport feature could also affect the validity of SYN cookies (see Section 3.6 of [RFC4987]). Implementation over UDP: do nothing.

- o Set maximum burst Protocols: SCTP Automatable because it relates to knowledge about the network, not the application.
- o Configure size where messages are broken up for partial delivery Protocols: SCTP Functional because this is closely tied to properties of the data that an application sends or expects to receive. Implementation over TCP: not possible. Implementation over UDP: not possible.
- o Disable checksum when sending Protocols: UDP Functional because application-specific knowledge is necessary to decide whether it can be acceptable to lose data integrity. Implementation: via SET_CHECKSUM_ENABLED.UDP. Implementation over TCP: do nothing.
- o Disable checksum requirement when receiving Protocols: UDP Functional because application-specific knowledge is necessary to decide whether it can be acceptable to lose data integrity. Implementation: via SET_CHECKSUM_REQUIRED.UDP. Implementation over TCP: do nothing.
- o Specify checksum coverage used by the sender Protocols: UDP-Lite

Functional because application-specific knowledge is necessary to decide for which parts of the data it can be acceptable to lose data integrity. Implementation: via SET_CHECKSUM_COVERAGE.UDP-Lite. Implementation over TCP: do nothing.

- o Specify minimum checksum coverage required by receiver Protocols: UDP-Lite Functional because application-specific knowledge is necessary to decide for which parts of the data it can be acceptable to lose data integrity. Implementation: via SET_MIN_CHECKSUM_COVERAGE.UDP-Lite. Implementation over TCP: do nothing.
- o Specify DF field Protocols: UDP(-Lite) Optimizing because the DF field can be used to carry out Path MTU Discovery, which can lead an application to choose message sizes that can be transmitted more efficiently. Implementation: via MAINTENANCE.SET_DF.UDP(-Lite) and SEND_FAILURE.UDP(-Lite). Implementation over TCP: do nothing. With TCP the sender is not in control of transport message sizes, making this functionality irrelevant.
- o Get max. transport-message size that may be sent using a nonfragmented IP packet from the configured interface Protocols: UDP(-Lite) Optimizing because this can lead an application to choose message sizes that can be transmitted more efficiently. Implementation over TCP: do nothing: this information is not available with TCP.
- o Get max. transport-message size that may be received from the configured interface Protocols: UDP(-Lite) Optimizing because this can, for example, influence an application's memory management. Implementation over TCP: do nothing: this information is not available with TCP.

o Specify TTL/Hop count field Protocols: UDP(-Lite) Automatable because a transport system can use a large enough system default to avoid communication failures. Allowing an application to configure it differently can produce notifications of ICMP error message arrivals that yield information which only relates to knowledge about the network, not the application.

- o Obtain TTL/Hop count field Protocols: UDP(-Lite) Automatable because the TTL/Hop count field relates to knowledge about the network, not the application.
- o Specify ECN field Protocols: UDP(-Lite) Automatable because the ECN field relates to knowledge about the network, not the application.
- o Obtain ECN field Protocols: UDP(-Lite) Optimizing because this information can be used by an application to better carry out congestion control (this is relevant when choosing a data transmission transport service that does not already do congestion control). Implementation over TCP: do nothing: this information is not available with TCP.
- o Specify IP Options Protocols: UDP(-Lite) Automatable because IP Options relate to knowledge about the network, not the application.
- o Obtain IP Options Protocols: UDP(-Lite) Automatable because IP Options relate to knowledge about the network, not the application.

o Enable and configure a "Low Extra Delay Background Transfer" Protocols: A protocol implementing the LEDBAT congestion control mechanism Optimizing because whether this service is appropriate or not depends on application-specific knowledge. However, wrongly using this will only affect the speed of data transfers (albeit including other transfers that may compete with the transport system's transfer in the network), so it is still correct within the "best effort" service model. Implementation: via CONFIGURE.LEDBAT and/or SET_DSCP.TCP / SET_DSCP.SCTP / SET_DSCP.UDP(-Lite) [LBE-draft]. Implementation over TCP: do nothing. Implementation over UDP: do nothing.

TERMINATION:

- o Close after reliably delivering all remaining data, causing an event informing the application on the other side Protocols: TCP, SCTP Functional because the notion of a connection is often reflected in applications as an expectation to have all outstanding data delivered and no longer be able to communicate after a "Close" succeeded, with a communication sequence relating to this transport feature that is defined by the application protocol. Implementation: via CLOSE.TCP and CLOSE.SCTP. Implementation over UDP: not possible.
- o Abort without delivering remaining data, causing an event informing the application on the other side Protocols: TCP, SCTP Functional because the notion of a connection is often reflected in applications as an expectation to potentially not have all outstanding data delivered and no longer be able to communicate after an "Abort" succeeded. On both sides of a connection, an application protocol may define a communication sequence relating to this transport feature. Implementation: via ABORT.TCP and ABORT.SCTP. Implementation over UDP: not possible.
- o Abort without delivering remaining data, not causing an event informing the application on the other side

Protocols: UDP(-Lite) Functional because the notion of a connection is often reflected in applications as an expectation to potentially not have all outstanding data delivered and no longer be able to communicate after an "Abort" succeeded. On both sides of a connection, an application protocol may define a communication sequence relating to this transport feature. Implementation: via ABORT.UDP(-Lite). Implementation over TCP: stop using the connection, wait for a timeout.

o Timeout event when data could not be delivered for too long Protocols: TCP, SCTP Functional because this notifies that potentially assumed reliable data delivery is no longer provided. Implementation: via TIMEOUT.TCP and TIMEOUT.SCTP. Implementation over UDP: do nothing: this event will not occur with UDP.

A.1.2. DATA Transfer Related Transport Features

A.1.2.1. Sending Data

o Reliably transfer data, with congestion control Protocols: TCP, SCTP Functional because this is closely tied to properties of the data that an application sends or expects to receive. Implementation: via SEND.TCP and SEND.SCTP. Implementation over UDP: not possible.

o Reliably transfer a message, with congestion control Protocols: SCTP Functional because this is closely tied to properties of the data that an application sends or expects to receive. Implementation: via SEND.SCTP. Implementation over TCP: via SEND.TCP. With SEND.TCP, messages will not be identifiable by the receiver. Implementation over UDP: not possible.

Internet-Draft

- o Unreliably transfer a message Protocols: SCTP, UDP(-Lite) Optimizing because only applications know about the time criticality of their communication, and reliably transfering a message is never incorrect for the receiver of a potentially unreliable data transfer, it is just slower. ADDED. This differs from the 2 automatable transport features below in that it leaves the choice of congestion control open. Implementation: via SEND.SCTP or SEND.UDP(-Lite). Implementation over TCP: use SEND.TCP. With SEND.TCP, messages will be sent reliably, and they will not be identifiable by the receiver.
- o Unreliably transfer a message, with congestion control Protocols: SCTP Automatable because congestion control relates to knowledge about the network, not the application.
- o Unreliably transfer a message, without congestion control Protocols: UDP(-Lite) Automatable because congestion control relates to knowledge about the network, not the application.
- o Configurable Message Reliability Protocols: SCTP Optimizing because only applications know about the time criticality of their communication, and reliably transfering a message is never incorrect for the receiver of a potentially unreliable data transfer, it is just slower. Implementation: via SEND.SCTP. Implementation over TCP: By using SEND.TCP and ignoring this configuration: based on the assumption of the best-effort service model, unnecessarily delivering data does not violate application expectations. Moreover, it is not possible to associate the requested reliability to a "message" in TCP anyway. Implementation over UDP: not possible.

o Choice of stream Protocols: SCTP

Automatable because it requires using multiple streams, but requesting multiple streams in the CONNECTION.ESTABLISHMENT category is automatable. Implementation: see Appendix A.3.2.

- o Choice of path (destination address) Protocols: SCTP Automatable because it requires using multiple sockets, but obtaining multiple sockets in the CONNECTION.ESTABLISHMENT category is automatable.
- o Ordered message delivery (potentially slower than unordered) Protocols: SCTP Functional because this is closely tied to properties of the data that an application sends or expects to receive. Implementation: via SEND.SCTP. Implementation over TCP: By using SEND.TCP. With SEND.TCP, messages will not be identifiable by the receiver. Implementation over UDP: not possible.
- o Unordered message delivery (potentially faster than ordered) Protocols: SCTP, UDP(-Lite) Functional because this is closely tied to properties of the data that an application sends or expects to receive. Implementation: via SEND.SCTP. Implementation over TCP: By using SEND.TCP and always sending data ordered: based on the assumption of the best-effort service model, ordered delivery may just be slower and does not violate application expectations. Moreover, it is not possible to associate the requested delivery order to a "message" in TCP anyway.
- o Request not to bundle messages Protocols: SCTP Optimizing because this decision depends on knowledge about the size of future data blocks and the delay between them. Implementation: via SEND.SCTP. Implementation over TCP: By using SEND.TCP and DISABLE_NAGLE.TCP to disable the Nagle algorithm when the request is made and enable it again when the request is no longer made. Note that this is not fully equivalent because it relates to the time of issuing the request rather than a specific message.

Implementation over UDP: do nothing (UDP never bundles messages).

o Specifying a "payload protocol-id" (handed over as such by the receiver) Protocols: SCTP Functional because it allows to send extra application data with every message, for the sake of identification of data, which by itself is application-specific. Implementation: SEND.SCTP. Implementation over TCP: not possible. Implementation over UDP: not possible.

- o Specifying a key id to be used to authenticate a message Protocols: SCTP Functional because this has a direct influence on security. Implementation: via a parameter in SEND.SCTP. Implementation over TCP: This could be emulated by using SET_AUTH.TCP before and after the message is sent. Note that this is not fully equivalent because it relates to the time of issuing the request rather than a specific message. Implementation over UDP: not possible.
- o Request not to delay the acknowledgement (SACK) of a message Protocols: SCTP Optimizing because only an application knows for which message it wants to quickly be informed about success / failure of its delivery. Implementation over TCP: do nothing. Implementation over UDP: do nothing.

A.1.2.2. Receiving Data

o Receive data (with no message delimiting) Protocols: TCP Functional because a transport system must be able to send and receive data. Implementation: via RECEIVE.TCP. Implementation over UDP: do nothing (hand over a message, let the application ignore message boundaries).

o Receive a message Protocols: SCTP, UDP(-Lite) Functional because this is closely tied to properties of the data that an application sends or expects to receive. Implementation: via RECEIVE.SCTP and RECEIVE.UDP(-Lite). Implementation over TCP: not possible.

o Choice of stream to receive from Protocols: SCTP Automatable because it requires using multiple streams, but requesting multiple streams in the CONNECTION.ESTABLISHMENT category is automatable. Implementation: see Appendix A.3.2.

o Information about partial message arrival Protocols: SCTP Functional because this is closely tied to properties of the data that an application sends or expects to receive. Implementation: via RECEIVE.SCTP. Implementation over TCP: do nothing: this information is not available with TCP. Implementation over UDP: do nothing: this information is not available with UDP.

A.1.2.3. Errors

This section describes sending failures that are associated with a specific call to in the "Sending Data" category (Appendix A.1.2.1).

o Notification of send failures Protocols: SCTP, UDP(-Lite) Functional because this notifies that potentially assumed reliable data delivery is no longer provided. ADDED. This differs from the 2 automatable transport features below in that it does not distinugish between unsent and unacknowledged messages. Implementation: via SENDFAILURE-EVENT.SCTP and SEND_FAILURE.UDP(-Lite). Implementation over TCP: do nothing: this notification is not available and will therefore not occur with TCP.

- Notification of an unsent (part of a) message
 Protocols: SCTP, UDP(-Lite)
 Automatable because the distinction between unsent and unacknowledged is network-specific.
- Notification of an unacknowledged (part of a) message Protocols: SCTP
 Automatable because the distinction between unsent and unacknowledged is network-specific.
- Notification that the stack has no more user data to send Protocols: SCTP
 Optimizing because reacting to this notification requires the application to be involved, and ensuring that the stack does not run dry of data (for too long) can improve performance.
 Implementation over TCP: do nothing. See also the discussion in Appendix A.3.4.
 Implementation over UDP: do nothing. This notification is not available and will therefore not occur with UDP.
- Notification to a receiver that a partial message delivery has been aborted Protocols: SCTP Functional because this is closely tied to properties of the data that an application sends or expects to receive. Implementation over TCP: do nothing. This notification is not available and will therefore not occur with TCP. Implementation over UDP: do nothing. This notification is not available and will therefore not occur with UDP.

A.2. Step 2: Reduction -- The Reduced Set of Transport Features

By hiding automatable transport features from the application, a transport system can gain opportunities to automate the usage of network-related functionality. This can facilitate using the transport system for the application programmer and it allows for optimizations that may not be possible for an application. For instance, system-wide configurations regarding the usage of multiple interfaces can better be exploited if the choice of the interface is
Minimal TAPS Transport Services February 2018

not entirely up to the application. Therefore, since they are not strictly necessary to expose in a transport system, we do not include automatable transport features in the reduced set of transport features. This leaves us with only the transport features that are either optimizing or functional.

A transport system should be able to communicate via TCP or UDP if alternative transport protocols are found not to work. For many transport features, this is possible -- often by simply not doing anything when a specific request is made. For some transport features, however, it was identified that direct usage of neither TCP nor UDP is possible: in these cases, even not doing anything would incur semantically incorrect behavior. Whenever an application would make use of one of these transport features, this would eliminate the possibility to use TCP or UDP. Thus, we only keep the functional and optimizing transport features for which an implementation over either TCP or UDP is possible in our reduced set.

In the following list, we precede a transport feature with "T:" if an implementation over TCP is possible, "U:" if an implementation over UDP is possible, and "TU:" if an implementation over either TCP or UDP is possible.

A.2.1. CONNECTION Related Transport Features

ESTABLISHMENT:

- o T,U: Connect
- o T,U: Specify number of attempts and/or timeout for the first establishment message
- o T: Configure authentication
- o T: Hand over a message to reliably transfer (possibly multiple times) before connection establishment
- o T: Hand over a message to reliably transfer during connection establishment

AVAILABILITY:

- o T,U: Listen
- o T: Configure authentication

MAINTENANCE:

- o T: Change timeout for aborting connection (using retransmit limit or time value)
- o T: Suggest timeout to the peer
- o T,U: Disable Nagle algorithm

- o T,U: Notification of Excessive Retransmissions (early warning below abortion threshold)
- o T,U: Specify DSCP field
- o T,U: Notification of ICMP error message arrival
- o T: Change authentication parameters
- o T: Obtain authentication information
- o T,U: Set Cookie life value
- o T,U: Choose a scheduler to operate between streams of an association
- o T,U: Configure priority or weight for a scheduler
- o T,U: Disable checksum when sending
- o T,U: Disable checksum requirement when receiving
- o T,U: Specify checksum coverage used by the sender
- o T,U: Specify minimum checksum coverage required by receiver
- o T,U: Specify DF field
- o T,U: Get max. transport-message size that may be sent using a nonfragmented IP packet from the configured interface
- o T,U: Get max. transport-message size that may be received from the configured interface
- o T,U: Obtain ECN field
- o T,U: Enable and configure a "Low Extra Delay Background Transfer"

TERMINATION:

- o T: Close after reliably delivering all remaining data, causing an event informing the application on the other side
- o T: Abort without delivering remaining data, causing an event informing the application on the other side
- o T,U: Abort without delivering remaining data, not causing an event informing the application on the other side
- o T,U: Timeout event when data could not be delivered for too long

A.2.2. DATA Transfer Related Transport Features

A.2.2.1. Sending Data

- o T: Reliably transfer data, with congestion control
- o T: Reliably transfer a message, with congestion control
- o T,U: Unreliably transfer a message
- o T: Configurable Message Reliability
- o T: Ordered message delivery (potentially slower than unordered)
- o T,U: Unordered message delivery (potentially faster than ordered)
- o T,U: Request not to bundle messages
- o T: Specifying a key id to be used to authenticate a message
- o T,U: Request not to delay the acknowledgement (SACK) of a message

A.2.2.2. Receiving Data

- o T,U: Receive data (with no message delimiting)
- o U: Receive a message
- o T,U: Information about partial message arrival

A.2.2.3. Errors

This section describes sending failures that are associated with a specific call to in the "Sending Data" category (Appendix A.1.2.1).

- o T,U: Notification of send failures
- o T,U: Notification that the stack has no more user data to send
- o T,U: Notification to a receiver that a partial message delivery has been aborted

A.3. Step 3: Discussion

The reduced set in the previous section exhibits a number of peculiarities, which we will discuss in the following. This section focuses on TCP because, with the exception of one particular transport feature ("Receive a message" -- we will discuss this in <u>Appendix A.3.1</u>), the list shows that UDP is strictly a subset of TCP. We can first try to understand how to build a transport system that can run over TCP, and then narrow down the result further to allow that the system can always run over either TCP or UDP (which effectively means removing everything related to reliability, ordering, authentication and closing/aborting with a notification to the peer).

Note that, because the functional transport features of UDP are -with the exception of "Receive a message" -- a subset of TCP, TCP can be used as a replacement for UDP whenever an application does not need message delimiting (e.g., because the application-layer protocol already does it). This has been recognized by many applications that already do this in practice, by trying to communicate with UDP at first, and falling back to TCP in case of a connection failure.

A.3.1. Sending Messages, Receiving Bytes

For implementing a transport system over TCP, there are several transport features related to sending, but only a single transport feature related to receiving: "Receive data (with no message delimiting)" (and, strangely, "information about partial message arrival"). Notably, the transport feature "Receive a message" is also the only non-automatable transport feature of UDP(-Lite) for which no implementation over TCP is possible.

To support these TCP receiver semantics, we define an "Application-Framed Bytestream" (AFra-Bytestream). AFra-Bytestreams allow senders to operate on messages while minimizing changes to the TCP socket API. In particular, nothing changes on the receiver side - data can be accepted via a normal TCP socket.

In an AFra-Bytestream, the sending application can optionally inform the transport about message boundaries and required properties per message (configurable order and reliability, or embedding a request not to delay the acknowledgement of a message). Whenever the sending application specifies per-message properties that relax the notion of reliable in-order delivery of bytes, it must assume that the receiving application is 1) able to determine message boundaries, provided that messages are always kept intact, and 2) able to accept these relaxed per-message properties. Any signaling of such information to the peer is up to an application-layer protocol and considered out of scope of this document.

For example, if an application requests to transfer fixed-size messages of 100 bytes with partial reliability, this needs the receiving application to be prepared to accept data in chunks of 100 bytes. If, then, some of these 100-byte messages are missing (e.g., if SCTP with Configurable Reliability is used), this is the expected application behavior. With TCP, no messages would be missing, but this is also correct for the application, and the possible retransmission delay is acceptable within the best effort service model [<u>RFC7305</u>]. Still, the receiving application would separate the byte stream into 100-byte chunks.

Note that this usage of messages does not require all messages to be equal in size. Many application protocols use some form of Type-Length-Value (TLV) encoding, e.g. by defining a header including length fields; another alternative is the use of byte stuffing methods such as COBS [COBS]. If an application needs message numbers, e.g. to restore the correct sequence of messages, these must also be encoded by the application itself, as the sequence number related transport features of SCTP are not provided by the "minimum set" (in the interest of enabling usage of TCP).

A.3.2. Stream Schedulers Without Streams

We have already stated that multi-streaming does not require application-specific knowledge. Potential benefits or disadvantages of, e.g., using two streams of an SCTP association versus using two separate SCTP associations or TCP connections are related to knowledge about the network and the particular transport protocol in use, not the application. However, the transport features "Choose a scheduler to operate between streams of an association" and

"Configure priority or weight for a scheduler" operate on streams. Here, streams identify communication channels between which a scheduler operates, and they can be assigned a priority. Moreover, the transport features in the MAINTENANCE category all operate on assocations in case of SCTP, i.e. they apply to all streams in that assocation.

With only these semantics necessary to represent, the interface to a transport system becomes easier if we assume that connections may be a transport protocol's connection or association, but could also be a stream of an existing SCTP association, for example. We only need to allow for a way to define a possible grouping of connections. Then, all MAINTENANCE transport features can be said to operate on connection groups, not connections, and a scheduler operates on the connections within a group.

To be compatible with multiple transport protocols and uniformly allow access to both transport connections and streams of a multistreaming protocol, the semantics of opening and closing need to be the most restrictive subset of all of the underlying options. For example, TCP's support of half-closed connections can be seen as a feature on top of the more restrictive "ABORT"; this feature cannot be supported because not all protocols used by a transport system (including streams of an association) support half-closed connections.

A.3.3. Early Data Transmission

There are two transport features related to transferring a message early: "Hand over a message to reliably transfer (possibly multiple times) before connection establishment", which relates to TCP Fast Open [RFC7413], and "Hand over a message to reliably transfer during connection establishment", which relates to SCTP's ability to transfer data together with the COOKIE-Echo chunk. Also without TCP Fast Open, TCP can transfer data during the handshake, together with the SYN packet -- however, the receiver of this data may not hand it over to the application until the handshake has completed. Also, different from TCP Fast Open, this data is not delimited as a message by TCP (thus, not visible as a ``message''). This functionality is commonly available in TCP and supported in several implementations, even though the TCP specification does not explain how to provide it to applications.

A transport system could differentiate between the cases of transmitting data "before" (possibly multiple times) or "during" the handshake. Alternatively, it could also assume that data that are handed over early will be transmitted as early as possible, and "before" the handshake would only be used for messages that are

explicitly marked as "idempotent" (i.e., it would be acceptable to transfer them multiple times).

The amount of data that can successfully be transmitted before or during the handshake depends on various factors: the transport protocol, the use of header options, the choice of IPv4 and IPv6 and the Path MTU. A transport system should therefore allow a sending application to query the maximum amount of data it can possibly transmit before (or, if exposed, during) connection establishment.

A.3.4. Sender Running Dry

The transport feature "Notification that the stack has no more user data to send" relates to SCTP's "SENDER DRY" notification. Such notifications can, in principle, be used to avoid having an unnecessarily large send buffer, yet ensure that the transport sender always has data available when it has an opportunity to transmit it. This has been found to be very beneficial for some applications [WWDC2015]. However, "SENDER DRY" truly means that the entire send buffer (including both unsent and unacknowledged data) has emptied -i.e., when it notifies the sender, it is already too late, the transport protocol already missed an opportunity to send data. Some modern TCP implementations now include the unspecified "TCP_NOTSENT_LOWAT" socket option that was proposed in [WWDC2015], which limits the amount of unsent data that TCP can keep in the socket buffer; this allows to specify at which buffer filling level the socket becomes writable, rather than waiting for the buffer to run empty.

SCTP allows to configure the sender-side buffer too: the automatable Transport Feature "Configure send buffer size" provides this functionality, but only for the complete buffer, which includes both unsent and unacknowledged data. SCTP does not allow to control these two sizes separately. It therefore makes sense for a transport system to allow for uniform access to "TCP_NOTSENT_LOWAT" as well as the "SENDER DRY" notification.

A.3.5. Capacity Profile

The transport features:

- o Disable Nagle algorithm
- o Enable and configure a "Low Extra Delay Background Transfer"
- o Specify DSCP field

all relate to a QoS-like application need such as "low latency" or "scavenger". In the interest of flexibility of a transport system, they could therefore be offered in a uniform, more abstract way,

where a transport system could e.g. decide by itself how to use combinations of LEDBAT-like congestion control and certain DSCP values, and an application would only specify a general "capacity profile" (a description of how it wants to use the available capacity). A need for "lowest possible latency at the expense of overhead" could then translate into automatically disabling the Nagle algorithm.

In some cases, the Nagle algorithm is best controlled directly by the application because it is not only related to a general profile but also to knowledge about the size of future messages. For fine-grain control over Nagle-like functionality, the "Request not to bundle messages" is available.

A.3.6. Security

Both TCP and SCTP offer authentication. TCP authenticates complete segments. SCTP allows to configure which of SCTP's chunk types must always be authenticated -- if this is exposed as such, it creates an undesirable dependency on the transport protocol. For compatibility with TCP, a transport system should only allow to configure complete transport layer packets, including headers, IP pseudo-header (if any) and payload.

Security is discussed in a separate TAPS document [<u>I-D.pauly-taps-transport-security</u>]. The minimal set presented in the present document therefore excludes all security related transport features: "Configure authentication", "Change authentication parameters", "Obtain authentication information" and and "Set Cookie life value" as well as "Specifying a key id to be used to authenticate a message".

A.3.7. Packet Size

UDP(-Lite) has a transport feature called "Specify DF field". This yields an error message in case of sending a message that exceeds the Path MTU, which is necessary for a UDP-based application to be able to implement Path MTU Discovery (a function that UDP-based applications must do by themselves). The "Get max. transport-message size that may be sent using a non-fragmented IP packet from the configured interface" transport feature yields an upper limit for the Path MTU (minus headers) and can therefore help to implement Path MTU Discovery more efficiently.

Appendix B. Revision information

XXX RFC-Ed please remove this section prior to publication.

-02: implementation suggestions added, discussion section added, terminology extended, DELETED category removed, various other fixes; list of Transport Features adjusted to -01 version of [RFC8303] except that MPTCP is not included.

-03: updated to be consistent with -02 version of [RFC8303].

-04: updated to be consistent with -03 version of [RFC8303]. Reorganized document, rewrote intro and conclusion, and made a first stab at creating a real "minimal set".

-05: updated to be consistent with -05 version of [RFC8303] (minor changes). Fixed a mistake regarding Cookie Life value. Exclusion of security related transport features (to be covered in a separate document). Reorganized the document (now begins with the minset, derivation is in the appendix). First stab at an abstract API for the minset.

draft-ietf-taps-minset-00: updated to be consistent with -08 version of [RFC8303] ("obtain message delivery number" was removed, as this has also been removed in [RFC8303] because it was a mistake in RFC4960. This led to the removal of two more transport features that were only designated as functional because they affected "obtain message delivery number"). Fall-back to UDP incorporated (this was requested at IETF-99); this also affected the transport feature "Choice between unordered (potentially faster) or ordered delivery of messages" because this is a boolean which is always true for one fall-back protocol, and always false for the other one. This was therefore now divided into two features, one for ordered, one for unordered delivery. The word "reliably" was added to the transport features "Hand over a message to reliably transfer (possibly multiple times) before connection establishment" and "Hand over a message to reliably transfer during connection establishment" to make it clearer why this is not supported by UDP. Clarified that the "minset abstract interface" is not proposing a specific API for all TAPS systems to implement, but it is just a way to describe the minimum set. Author order changed.

WG -01: "fall-back to" (TCP or UDP) replaced (mostly with "implementation over"). References to post-sockets removed (these were statments that assumed that post-sockets requires two-sided implementation). Replaced "flow" with "TAPS Connection" and "frame" with "message" to avoid introducing new terminology. Made sections 3 and 4 in line with the categorization that is already used in the

appendix and [RFC8303], and changed style of section 4 to be even shorter and less interface-like. Updated reference draft-ietf-tsvwgsctp-ndata to RFC8260.

WG -02: rephrased "the TAPS system" and "TAPS connection" etc. to more generally talk about transport after the intro (mostly replacing "TAPS system" with "transport system" and "TAPS connection" with "connection". Merged sections $\underline{3}$ and $\underline{4}$ to form a new section $\underline{3}$.

Authors' Addresses

Michael Welzl University of Oslo PO Box 1080 Blindern Oslo N-0316 Norway

Phone: +47 22 85 24 20 Email: michawe@ifi.uio.no

Stein Gjessing University of Oslo PO Box 1080 Blindern Oslo N-0316 Norway

Phone: +47 22 85 24 44 Email: steing@ifi.uio.no