

taps
Internet-Draft
Intended status: Informational
Expires: 28 April 2022

M. Duke
F5 Networks, Inc.
25 October 2021

TAPS Transport Discovery
draft-duke-taps-transport-discovery-02

Abstract

The Transport Services architecture decouples applications from the protocol implementations that transport their data. While it is often straightforward to connect applications with transports that are present in the host operating system, providing a means of discovering user-installed implementations dramatically enlarges the use cases. This document discusses considerations for the design of a discovery mechanism and an example of such a design.

Discussion of this work is encouraged to happen on the TAPS IETF mailing list taps@ietf.org or on the GitHub repository which contains the draft: <https://github.com/martinduke/draft-duke-taps-transport-discovery>.

Discussion Venues

This note is to be removed before publishing as an RFC.

Discussion of this document takes place on the mailing list (taps@ietf.org), which is archived at <https://mailarchive.ietf.org/arch/browse/taps/>.

Source for this draft and an issue tracker can be found at <https://github.com/martinduke/draft-duke-taps-transport-discovery>.

Status of This Memo

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

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 <https://datatracker.ietf.org/drafts/current/>.

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 28 April 2022.

Copyright Notice

Copyright (c) 2021 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 (<https://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.

Table of Contents

1. Introduction	3
2. Conventions	4
3. Entities	4
4. Protocol Implementation	5
4.1. Functions	5
4.2. Events	5
5. Protocol Installer	6
6. TAPS	6
7. Security Considerations	7
8. IANA Considerations	8
9. Implementation Status	8
10. Informative References	8
Appendix A. Acknowledgments	8
Appendix B. Change Log	9
B.1. since draft-duke-taps-transport-discovery-01	9
Author's Address	9

1. Introduction

The Transport Services architecture [I-D.ietf-taps-arch] enables applications to be protocol-agnostic by presenting an interface where applications can specify their required properties, and the service will select whichever protocol implementation available in the system best meets those requirements. This increases application portability and eases the introduction of new transport innovations by not requiring changes to applications.

It is sometimes straightforward for a Transport Services interface to identify the transports available in the host operating system. However, including transports installed by the user greatly expands use cases for the architecture. This document presents considerations for the secure design of a system for discovery of new protocol implementations.

Protocol Discovery would ideally have several desirable properties.

- * The transport services API should not have to recompile when installing new implementations. This would not only disrupt ongoing connections, but also involve ordinary users in the complex business of downloading and building source code.
- * It should support user-space implementations. Most protocol innovation begins with user space implementations, and many transports (e.g. TLS, HTTP, QUIC) are usually implemented outside the kernel long after reaching maturity.
- * Protocol Discovery should not subject ordinary users to security vulnerabilities. A new protocol installation is an opportunity to hijack a user's networking stack, and Protocol Discovery requires strong protections against arbitrary code performing operations other than advertised on application data.
- * Conversely, sophisticated users need a means of discovering implementations that are too new to have fully developed internet trust mechanisms. This is the only means of initially deploying new protocols for existing apps, and is the most plausible model to deploy transport services API shims for existing protocol libraries (e.g., the common TLS implementations) before their proponents deploy native support.
- * Applications should not have to bring their own implementations. The Transport Services API has the concept of "framers" (see Sec. 7.1 of [I-D.ietf-taps-interface]) that provide some ability for applications to provide additional protocol encapsulation around their messages. However, one important advantage of Transport

Services is that applications do not have to rely on a third-party implementation that might not offer long term support, or add to their footprint where a functionally equivalent protocol implementation is already present on the system.

This document attempts to resolve the tension between some of these properties.

2. Conventions

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in RFC 2119 [RFC2119].

"TAPS" is an abbreviation for the transport services API.

For brevity, this document will use "app" as a shorthand for "application."

As in other TAPS documents, the concept of a "transport protocol" is expanded beyond the traditional "transport layer" to include other protocols that encapsulate application data, such as TLS, HTTP, and Websockets.

3. Entities

The Transport Services API (TAPS) is responsible for matching protocol capabilities with application requirements, and mediating further app communication with the selected protocol implementation. In this document, it actively discovers what implementations are available in the system.

The protocol implementation instantiates the transport. In this document, it offers a dynamically linked library that conforms to standard interfaces so that TAPS can interchangeably interact with it. In practice, this may be a shim layer if the underlying implementation does not support TAPS.

The protocol installer, aside from installing the implementation library and/or a TAPS shim layer, also is responsible for notifying TAPS that the implementation is present, and what its capabilities are.

Finally, the application leverages TAPS to initiate, manage, and terminate communications with other endpoints. This document does not require any changes to application behavior beyond those in the core TAPS design.

More detailed requirements for each of these entities is below.

4. Protocol Implementation

The protocol implementation must offer a dynamically linked library that offers certain APIs. TAPS SHOULD, in its documentation, provide a template for the format of these functions.

4.1. Functions

The objects below need not follow the semantics of the TAPS application API. In particular, a "message" is unlikely to have all the property information described there, instead being a more primitive buffer in which raw data is stored.

```
''' Listener := Listen(localEndpoint) '''
```

Listen opens a socket and listens on the specified address, and returns a handle to the resulting listener.

```
''' Listener.Stop() '''
```

Stop causes the listener to stop accepting connections. Subsequent events will return handles to the resulting connection.

```
''' Connection.Send(Message) Connection.Receive(Message) '''
```

TAPS will provide a Message object for the protocol to either send, or use to store incoming data.

Further APIs are TBD.

4.2. Events

The protocol needs to throw all the events described in the TAPS Application API, although the return values may not exactly conform to the same semantics.

TAPS SHOULD provide an event framework that frees the protocol implementation from running its own thread for a polling loop. TAPS also SHOULD account for the possibility that the implementation may have its own polling architecture. If true, the protocol MUST conform to the API by translating its events into the signals or callbacks that TAPS expects.

5. Protocol Installer

The installer might use the operating system's package manager or "app store", or be a simple script. Besides installing the implementation, the installer also writes data to a registry that TAPS will access to discover the implementation.

This data will include:

- * the name of the supported protocol(s);
- * optionally, the versions of those protocols;
- * the path to the implementations TAPS-compliant library;
- * the properties that the protocol implementation supports, as described in Section 4.2 of [I-D.ietf-taps-interface]; and
- * information to authenticate the entry (see Section 7).

Of course, a de-installer should remove the appropriate registry entry.

A TAPS implementation SHOULD provide a template for this registry information.

One potential instantiation of this would have protocol installers write a file to a directory that, in a specified markup language, described the information above.

6. TAPS

TAPS creates a registry for protocol implementations, which might be a database or a directory. To prevent inadvertent security vulnerabilities, the host system SHOULD, at minimum, require administrative privileges to write to the registry.

No later than upon receipt of request for a Preconnection, TAPS MUST access the registry to determine the available protocols and their properties. It is perfectly valid for there to be multiple implementations of a protocol.

TAPS SHOULD validate entries in the registry using the provided authentication data.

One potential instantiation would start daemon that monitored the status of the registry. Upon any change to the registry, the daemon might:

- * authenticate any new entry in accordance with security policy;
- * verify that the required function handles are present;
- * run tests to verify the installation's claimed properties;
- * inform the user of the new protocol, requesting permission to trust it; and
- * write the information into shared memory for the use of Preconnections.

7. Security Considerations

User-space installation of protocols provides enormous opportunities for attackers to hijack a network stack. While this has always been possible with arbitrary protocol implementations, with TAPS applications completely unaware of the installation can be victims of such an attack.

An implementation might advertise properties it does not actually provide to attract more traffic. For example, a "TLS" implementation might not encrypt anything at all. A TAPS implementation MAY run tests on newly installed protocols to verify it provides the advertised properties.

Moreover, in principle an implementation could deliver application data anywhere it wanted with little visibility to the application, much less the user.

The origin of the protocol installer is important to the trust model. Obviously, transports in the kernel do not introduce vulnerabilities specific to TAPS. A trusted package manager (e.g. the Apple App Store or yum) may imply a minimal level of veracity of the available packages. Protocol implementations directly downloaded from the internet without mediation through these mechanisms require the greatest care.

Ongoing work on this document will largely focus on building mechanisms to mitigate this weakness. Some promising approaches include:

- * administrative privileges to alter the TAPS registry;
- * a special certificate authority that provides an authentication of the implementation's explicit and implicit claims, as well as the integrity of the installed binary;

- * each installer generates a private key and provides the corresponding public key, so that only possessors of the private key can modify or delete the registry entry;
- * confirmation by a human, prominently warned of potential consequences, if the installation is not mediated through a trusted authority.

8. IANA Considerations

This document has no IANA requirements.

9. Implementation Status

RFC Editor's Note: Please remove this section prior to publication of a final version of this document.

The Dynamic TAPS project (<https://github.com/f5networks/dynamic-taps>) is a preliminary effort to implement the concepts in this document.

10. Informative References

[I-D.ietf-taps-arch]

Pauly, T., Trammell, B., Brunstrom, A., Fairhurst, G., Perkins, C., Tiesel, P. S., and C. A. Wood, "An Architecture for Transport Services", Work in Progress, Internet-Draft, draft-ietf-taps-arch-11, 12 July 2021, <<https://datatracker.ietf.org/doc/html/draft-ietf-taps-arch-11>>.

[I-D.ietf-taps-interface]

Trammell, B., Welzl, M., Enghardt, T., Fairhurst, G., Kuehlewind, M., Perkins, C., Tiesel, P. S., Wood, C. A., Pauly, T., and K. Rose, "An Abstract Application Layer Interface to Transport Services", Work in Progress, Internet-Draft, draft-ietf-taps-interface-13, 12 July 2021, <<https://datatracker.ietf.org/doc/html/draft-ietf-taps-interface-13>>.

[RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://doi.org/10.17487/RFC2119>>.

Appendix A. Acknowledgments

Tim Worsley contributed important ideas to this document.

Appendix B. Change Log

RFC Editor's Note: Please remove this section prior to publication of a final version of this document.

B.1. since draft-duke-taps-transport-discovery-01

- * Added output of initial implementation work

Author's Address

Martin Duke
F5 Networks, Inc.

Email: martin.h.duke@gmail.com

TAPS Working Group
Internet-Draft
Intended status: Standards Track
Expires: 7 July 2022

T. Pauly, Ed.
Apple Inc.
B. Trammell, Ed.
Google Switzerland GmbH
A. Brunstrom
Karlstad University
G. Fairhurst
University of Aberdeen
C. Perkins
University of Glasgow
3 January 2022

An Architecture for Transport Services
draft-ietf-taps-arch-12

Abstract

This document describes an architecture for exposing transport protocol features to applications for network communication, a Transport Services system. The Transport Services Application Programming Interface (API) is based on an asynchronous, event-driven interaction pattern. This API uses messages for representing data transfer to applications, and describes how implementations can use multiple IP addresses, multiple protocols, and multiple paths, and provide multiple application streams. This document further defines common terminology and concepts to be used in definitions of a Transport Service API and a Transport Services implementation.

Status of This Memo

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

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 <https://datatracker.ietf.org/drafts/current/>.

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 7 July 2022.

Copyright Notice

Copyright (c) 2022 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 (<https://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 Revised BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Revised BSD License.

Table of Contents

1. Introduction	3
1.1. Background	4
1.2. Overview	4
1.3. Specification of Requirements	5
2. API Model	5
2.1. Event-Driven API	7
2.2. Data Transfer Using Messages	8
2.3. Flexible Implementation	8
3. API and Implementation Requirements	9
3.1. Provide Common APIs for Common Features	10
3.2. Allow Access to Specialized Features	11
3.3. Select Equivalent Protocol Stacks	12
3.4. Maintain Interoperability	13
4. Transport Services Architecture and Concepts	13
4.1. Transport Services API Concepts	14
4.1.1. Endpoint Objects	16
4.1.2. Connections and Related Objects	16
4.1.3. Pre-Establishment	18
4.1.4. Establishment Actions	18
4.1.5. Data Transfer Objects and Actions	19
4.1.6. Event Handling	20
4.1.7. Termination Actions	21
4.1.8. Connection Groups	21
4.2. Transport Services Implementation	22
4.2.1. Candidate Gathering	23
4.2.2. Candidate Racing	23
4.2.3. Separating Connection Contexts	24
5. IANA Considerations	24
6. Security and Privacy Considerations	25
7. Acknowledgements	26
8. References	26
8.1. Normative References	26

8.2. Informative References	26
Authors' Addresses	28

1. Introduction

Many application programming interfaces (APIs) to perform transport networking have been deployed, perhaps the most widely known and imitated being the BSD Socket [POSIX] interface (Socket API). The naming of objects and functions across these APIs is not consistent, and varies depending on the protocol being used. For example, sending and receiving streams of data is conceptually the same for both an unencrypted Transmission Control Protocol (TCP) stream and operating on an encrypted Transport Layer Security (TLS) [RFC8446] stream over TCP, but applications cannot use the same socket send() and recv() calls on top of both kinds of connections. Similarly, terminology for the implementation of transport protocols varies based on the context of the protocols themselves: terms such as "flow", "stream", "message", and "connection" can take on many different meanings. This variety can lead to confusion when trying to understand the similarities and differences between protocols, and how applications can use them effectively.

The goal of the Transport Services architecture is to provide a flexible and reusable architecture that provides a common interface for transport protocols. As applications adopt this interface, they will benefit from a wide set of transport features that can evolve over time, and ensure that the system providing the interface can optimize its behavior based on the application requirements and network conditions, without requiring changes to the applications. This flexibility enables faster deployment of new features and protocols. It can also support applications by offering racing mechanisms (attempting multiple IP addresses, protocols, or network paths in parallel), which otherwise need to be implemented in each application separately (see Section 4.2.2).

This document was developed in parallel with the specification of the Transport Services API [I-D.ietf-taps-interface] and implementation guidelines [I-D.ietf-taps-impl]. Although following the Transport Services architecture does not require that all APIs and implementations are identical, a common minimal set of features represented in a consistent fashion will enable applications to be easily ported from one system to another.

1.1. Background

The Transport Services architecture is based on the survey of services provided by IETF transport protocols and congestion control mechanisms [RFC8095], and the distilled minimal set of the features offered by transport protocols [RFC8923]. These documents identified common features and patterns across all transport protocols developed thus far in the IETF.

Since transport security is an increasingly relevant aspect of using transport protocols on the Internet, this architecture also considers the impact of transport security protocols on the feature-set exposed by Transport Services [RFC8922].

One of the key insights to come from identifying the minimal set of features provided by transport protocols [RFC8923] was that features either require application interaction and guidance (referred to in that document as Functional or Optimizing Features), or else can be handled automatically by a system implementing Transport Services (referred to as Automatable Features). Among the identified Functional and Optimizing Features, some were common across all or nearly all transport protocols, while others could be seen as features that, if specified, would only be useful with a subset of protocols, but would not harm the functionality of other protocols. For example, some protocols can deliver messages faster for applications that do not require messages to arrive in the order in which they were sent. However, this functionality needs to be explicitly allowed by the application, since reordering messages would be undesirable in many cases.

1.2. Overview

This document describes the Transport Services architecture in three sections:

- * Section 2 describes how the API model of Transport Services architecture differs from traditional socket-based APIs. Specifically, it offers asynchronous event-driven interaction, the use of messages for data transfer, and the flexibility to use different transport protocols and paths without requiring major changes to the application.
- * Section 3 explains the fundamental requirements for a Transport Services system. These principles are intended to make sure that transport protocols can continue to be enhanced and evolve without requiring significant changes by application developers.

- * Section 4 presents a diagram showing the Transport Services architecture and defines the concepts that are used by both the API [I-D.ietf-taps-interface] and implementation guidelines [I-D.ietf-taps-impl]. The Preconnection allows applications to configure Connection Properties.
- * Section 4 also presents how an abstract Connection is used to select a transport protocol instance such as TCP, UDP, or another transport. The Connection represents an object that can be used to send and receive messages.

1.3. Specification of Requirements

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in BCP 14 [RFC2119] [RFC8174] when, and only when, they appear in all capitals, as shown here.

2. API Model

The traditional model of using sockets for networking can be represented as follows:

- * Applications create connections and transfer data using the Socket API.
- * The Socket API provides the interface to the implementations of TCP and UDP (typically implemented in the system's kernel).
- * TCP and UDP in the kernel send and receive data over the available network-layer interfaces.
- * Sockets are bound directly to transport-layer and network-layer addresses, obtained via a separate resolution step, usually performed by a system-provided stub resolver.

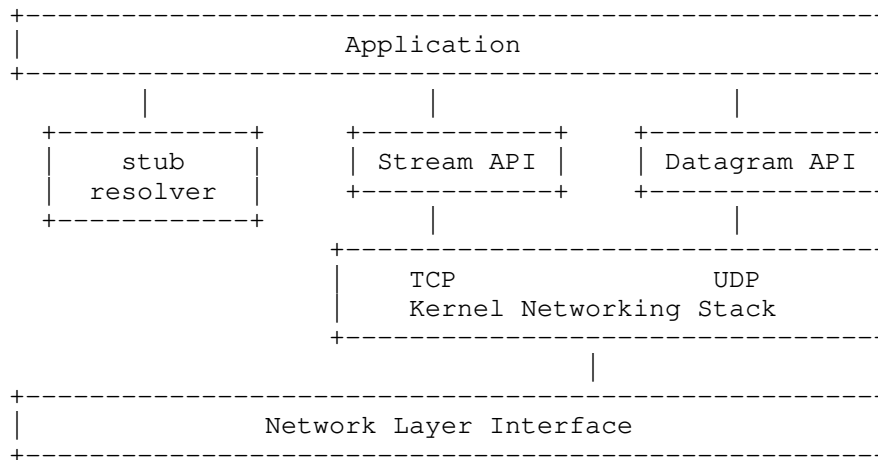


Figure 1: Socket API Model

The Transport Services architecture evolves this general model of interaction, to both modernize the API surface presented to applications by the transport layer and to enrich the capabilities of the implementation below the API.

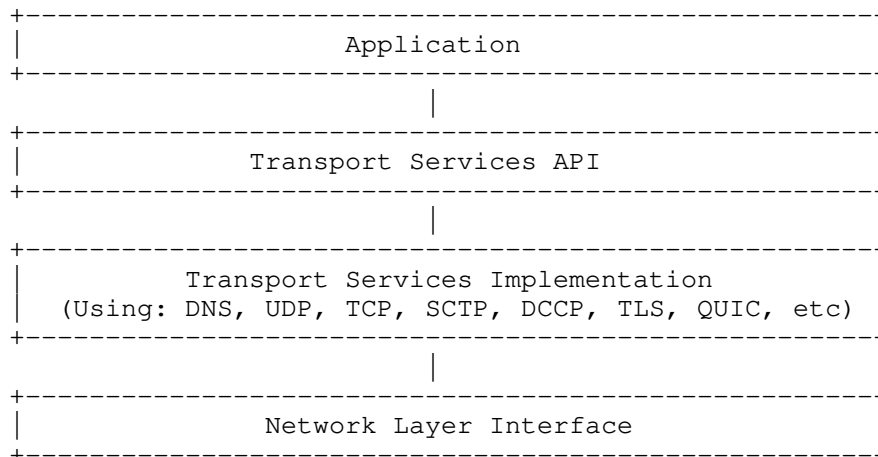


Figure 2: Transport Services API Model

The Transport Services API [I-D.ietf-taps-interface] defines the interface for an application to create Connections and transfer data. It combines interfaces for multiple interaction patterns into a unified whole. By combining name resolution with connection establishment and data transfer in a single API, it allows for more flexible implementations to provide path and transport protocol agility on the application's behalf.

The Transport Services implementation [I-D.ietf-taps-impl] implements the transport layer protocols and other functions needed to send and receive data. It is responsible for mapping the API to a specific available transport protocol stack and managing the available network interfaces and paths.

There are key differences between the Transport Services architecture and the architecture of the Socket API: the API of the Transport Services architecture is asynchronous and event-driven; it uses messages for representing data transfer to applications; and it describes how implementations can use multiple IP addresses, multiple protocols, multiple paths, and provide multiple application streams.

2.1. Event-Driven API

Originally, the Socket API presented a blocking interface for establishing connections and transferring data. However, most modern applications interact with the network asynchronously. Emulation of an asynchronous interface using the Socket API generally uses a try-and-fail model. If the application wants to read, but data has not yet been received from the peer, the call to read will fail. The application then waits and can try again later. In contrast to the Socket API, all interaction using the Transport Services API is expected to be asynchronous and use an event-driven model (see Section 4.1.6). For example, an application first issues a call to receive new data from the connection. When delivered data becomes available, this data is delivered to the application using asynchronous events that contain the data. Error handling is also asynchronous; a failure to send data results in an asynchronous error event.

This API also delivers events regarding the lifetime of a connection and changes in the available network links, which were not previously made explicit in the Socket API.

Using asynchronous events allows for a more natural interaction model when establishing connections and transferring data. Events in time more closely reflect the nature of interactions over networks, as opposed to how the Socket API represent network resources as file system objects that may be temporarily unavailable.

Separate from events, callbacks are also provided for asynchronous interactions with the API not directly related to events on the network or network interfaces.

2.2. Data Transfer Using Messages

The Socket API provides a message interface for datagram protocols like UDP, but provides an unstructured stream abstraction for TCP. While TCP has the ability to send and receive data as a byte-stream, most applications need to interpret structure within this byte-stream. For example, HTTP/1.1 uses character delimiters to segment messages over a byte-stream [RFC7230]; TLS record headers carry a version, content type, and length [RFC8446]; and HTTP/2 uses frames to segment its headers and bodies [RFC7540].

The Transport Services API represents data as messages, so that it more closely matches the way applications use the network. Providing a message-based abstraction provides many benefits, such as:

- * the ability to associate deadlines with messages, for applications that care about timing;
- * the ability control reliability, which messages to retransmit when there is packet loss, and how best to make use of the data that arrived;
- * the ability to automatically assign messages and connections to underlying transport connections to utilize multi-streaming and pooled connections.

Allowing applications to interact with messages is backwards-compatible with existing protocols and APIs because it does not change the wire format of any protocol. Instead, it gives the protocol stack additional information to allow it to make better use of modern transport services, while simplifying the application's role in parsing data. For protocols which natively use a streaming abstraction, framers (Section 4.1.5) bridge the gap between the two abstractions.

2.3. Flexible Implementation

The Socket API for protocols like TCP is generally limited to connecting to a single address over a single interface. It also presents a single stream to the application. Software layers built upon this API often propagate this limitation of a single-address single-stream model. The Transport Services architecture is designed:

- * to handle multiple candidate endpoints, protocols, and paths;
- * to support candidate protocol racing to select the most optimal stack in each situation;
- * to support multipath and multistreaming protocols;
- * to provide state caching and application control over it.

A Transport Services implementation is intended to be flexible at connection establishment time, considering many different options and trying to select the most optimal combinations by racing them and measuring the results (see Section 4.2.1 and Section 4.2.2). This requires applications to provide higher-level endpoints than IP addresses, such as hostnames and URLs, which are used by a Transport Services implementation for resolution, path selection, and racing. An implementation can further implement fallback mechanisms if connection establishment of one protocol fails or performance is detected to be unsatisfactory.

Information used in connection establishment (e.g. cryptographic resumption tokens, information about usability of certain protocols on the path, results of racing in previous connections) are cached in the Transport Services implementation. Applications have control over whether this information is used for a specific establishment, in order to allow tradeoffs between efficiency and linkability.

Flexibility after connection establishment is also important. Transport protocols that can migrate between multiple network-layer interfaces need to be able to process and react to interface changes. Protocols that support multiple application-layer streams need to support initiating and receiving new streams using existing connections.

3. API and Implementation Requirements

A goal of the Transport Services architecture is to redefine the interface between applications and transports in a way that allows the transport layer to evolve and improve without fundamentally changing the contract with the application. This requires a careful consideration of how to expose the capabilities of protocols. This architecture also encompasses system policies that can influence and inform how transport protocols use a network path or interface.

There are several ways the Transport Services system can offer flexibility to an application: it can provide access to transport protocols and protocol features; it can use these protocols across multiple paths that could have different performance and functional

characteristics; and it can communicate with different remote systems to optimize performance, robustness to failure, or some other metric. Beyond these, if the Transport Services API remains the same over time, new protocols and features can be added to the Transport Services implementation without requiring changes in applications for adoption. Similarly, this can provide a common basis for utilizing information about a network path or interface, enabling evolution below the transport layer.

The normative requirements described in this section allow Transport Services APIs and Transport Services implementation to provide this functionality without causing incompatibility or introducing security vulnerabilities.

3.1. Provide Common APIs for Common Features

Any functionality that is common across multiple transport protocols SHOULD be made accessible through a unified set of calls using the Transport Services API. As a baseline, any Transport Services API SHOULD allow access to the minimal set of features offered by transport protocols [RFC8923].

An application can specify constraints and preferences for the protocols, features, and network interfaces it will use via Properties. Properties are used by an application to declare its preferences for how the transport service should operate at each stage in the lifetime of a connection. Transport Properties are subdivided into Selection Properties, which specify which paths and protocol stacks can be used and are preferred by the application; Connection Properties, which inform decisions made during connection establishment and fine-tune the established connection; and Message Properties, set on individual Messages.

It is RECOMMENDED that the Transport Services API offers properties that are common to multiple transport protocols. This enables a Transport Services implementation to appropriately select between protocols that offer equivalent features. Similarly, it is RECOMMENDED that the Properties offered by the Transport Services API are applicable to a variety of network layer interfaces and paths, which permits racing of different network paths without affecting the applications using the API. Each is expected to have a default value.

It is RECOMMENDED that the default values for Properties are selected to ensure correctness for the widest set of applications, while providing the widest set of options for selection. For example, since both applications that require reliability and those that do not require reliability can function correctly when a protocol

provides reliability, reliability ought to be enabled by default. As another example, the default value for a Property regarding the selection of network interfaces ought to permit as many interfaces as possible.

Applications using the Transport Services API are REQUIRED to be robust to the automated selection provided by the Transport Services implementation. This automated selection is constrained by the properties and preferences expressed by the application and requires applications to explicitly set properties that define any necessary constraints on protocol, path, and interface selection.

3.2. Allow Access to Specialized Features

There are applications that will need to control fine-grained details of transport protocols to optimize their behavior and ensure compatibility with remote systems. It is therefore RECOMMENDED that the Transport Services API and the Transport Services implementation permit more specialized protocol features to be used.

A specialized feature could be needed by an application only when using a specific protocol, and not when using others. For example, if an application is using TCP, it could require control over the User Timeout Option for TCP; these options would not take effect for other transport protocols. In such cases, the API ought to expose the features in such a way that they take effect when a particular protocol is selected, but do not imply that only that protocol could be used. For example, if the API allows an application to specify a preference to use the User Timeout Option, communication would not fail when a protocol such as QUIC is selected.

Other specialized features, however, can be strictly required by an application and thus constrain the set of protocols that can be used. For example, if an application requires support for automatic handover or failover for a connection, only protocol stacks that provide this feature are eligible to be used, e.g., protocol stacks that include a multipath protocol or a protocol that supports connection migration. A Transport Services API needs to allow applications to define such requirements and constrain the options available to a Transport Services implementation. Since such options are not part of the core/common features, it will generally be simple for an application to modify its set of constraints and change the set of allowable protocol features without changing the core implementation.

3.3. Select Equivalent Protocol Stacks

A Transport Services implementation can select Protocol Stacks based on the Selection and Connection Properties communicated by the application, along with any security parameters. If two different Protocol Stacks can be safely swapped, or raced in parallel (see Section 4.2.2), then they are considered to be "equivalent". Equivalent Protocol Stacks are defined as stacks that can provide the same Transport Properties and interface expectations as requested by the application.

The following two examples show non-equivalent Protocol Stacks:

- * If the application requires preservation of message boundaries, a Protocol Stack that runs UDP as the top-level interface to the application is not equivalent to a Protocol Stack that runs TCP as the top-level interface. A UDP stack would allow an application to read out message boundaries based on datagrams sent from the remote system, whereas TCP does not preserve message boundaries on its own, but needs a framing protocol on top to determine message boundaries.
- * If the application specifies that it requires reliable transmission of data, then a Protocol Stack using UDP without any reliability layer on top would not be allowed to replace a Protocol Stack using TCP.

The following example shows equivalent Protocol Stacks:

- * If the application does not require reliable transmission of data, then a Protocol Stack that adds reliability could be regarded as an equivalent Protocol Stack as long as providing this would not conflict with any other application-requested properties.

To ensure that security protocols are not incorrectly swapped, a Transport Services implementation MUST only select Protocol Stacks that meet application requirements ([RFC8922]). A Transport Services implementation SHOULD only race Protocol Stacks where the transport security protocols within the stacks are identical. A Transport Services implementation MUST NOT automatically fall back from secure protocols to insecure protocols, or to weaker versions of secure protocols. A Transport Services implementation MAY allow applications to explicitly specify that fallback to a specific other version of a protocol \, e.g., to allow fallback to TLS 1.2 if TLS 1.3 is not available.

3.4. Maintain Interoperability

It is important to note that neither the Transport Services API [I-D.ietf-taps-interface] nor the guidelines for the Transport Service implementation [I-D.ietf-taps-impl] define new protocols or protocol capabilities that affect what is communicated across the network. A Transport Services system **MUST NOT** require that a peer on the other side of a connection uses the same API or implementation. A Transport Services implementation acting as a connection initiator is able to communicate with any existing endpoint that implements the transport protocol(s) and all the required properties selected. Similarly, a Transport Services implementation acting as a listener can receive connections for any protocol that is supported from an existing initiator that implements the protocol, independent of whether the initiator uses the Transport Services architecture or not.

A Transport Services system makes decisions that select protocols and interfaces. In normal use, a given version of a Transport Services system **SHOULD** result in consistent protocol and interface selection decisions for the same network conditions given the same set of Properties. This is intended to provide predictable outcomes to the application using the API.

4. Transport Services Architecture and Concepts

This section and the remainder of this document describe the architecture non-normatively. The concepts defined in this document are intended primarily for use in the documents and specifications that describe the Transport Services system. This includes the architecture, the Transport Services API and the associated Transport Services implementation. While the specific terminology can be used in some implementations, it is expected that there will remain a variety of terms used by running code.

The architecture divides the concepts for Transport Services system into two categories:

1. API concepts, which are intended to be exposed to applications; and
2. System-implementation concepts, which are intended to be internally used by a Transport Services implementation.

The following diagram summarizes the top-level concepts in the architecture and how they relate to one another.

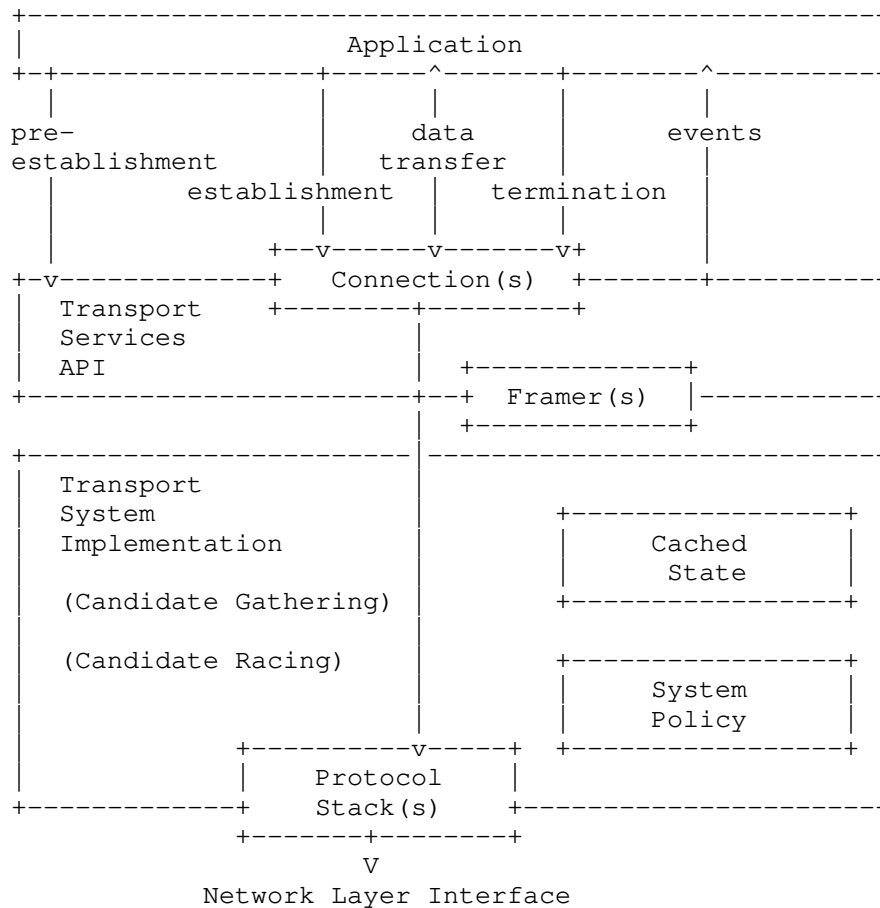


Figure 3: Concepts and Relationships in the Transport Services Architecture

4.1. Transport Services API Concepts

Fundamentally, a Transport Services API needs to provide connection objects (Section 4.1.2) that allow applications to establish communication, and then send and receive data. These could be exposed as handles or referenced objects, depending on the chosen programming language.

Beyond the connection objects, there are several high-level groups of actions that any Transport Services API needs to provide:

- * Pre-Establishment (Section 4.1.3) encompasses the properties that an application can pass to describe its intent, requirements, prohibitions, and preferences for its networking operations. These properties apply to multiple transport protocols, unless otherwise specified. Properties specified during Pre-Establishment can have a large impact on the rest of the interface: they modify how establishment occurs, they influence the expectations around data transfer, and they determine the set of events that will be supported.
- * Establishment (Section 4.1.4) focuses on the actions that an application takes on the connection objects to prepare for data transfer.
- * Data Transfer (Section 4.1.5) consists of how an application represents the data to be sent and received, the functions required to send and receive that data, and how the application is notified of the status of its data transfer.
- * Event Handling (Section 4.1.6) defines categories of notifications which an application can receive during the lifetime of transport objects. Events also provide opportunities for the application to interact with the underlying transport by querying state or updating maintenance options.
- * Termination (Section 4.1.7) focuses on the methods by which data transmission is stopped, and state is torn down in the transport.

The diagram below provides a high-level view of the actions and events during the lifetime of a Connection object. Note that some actions are alternatives (e.g., whether to initiate a connection or to listen for incoming connections), while others are optional (e.g., setting Connection and Message Properties in Pre-Establishment) or have been omitted for brevity and simplicity.

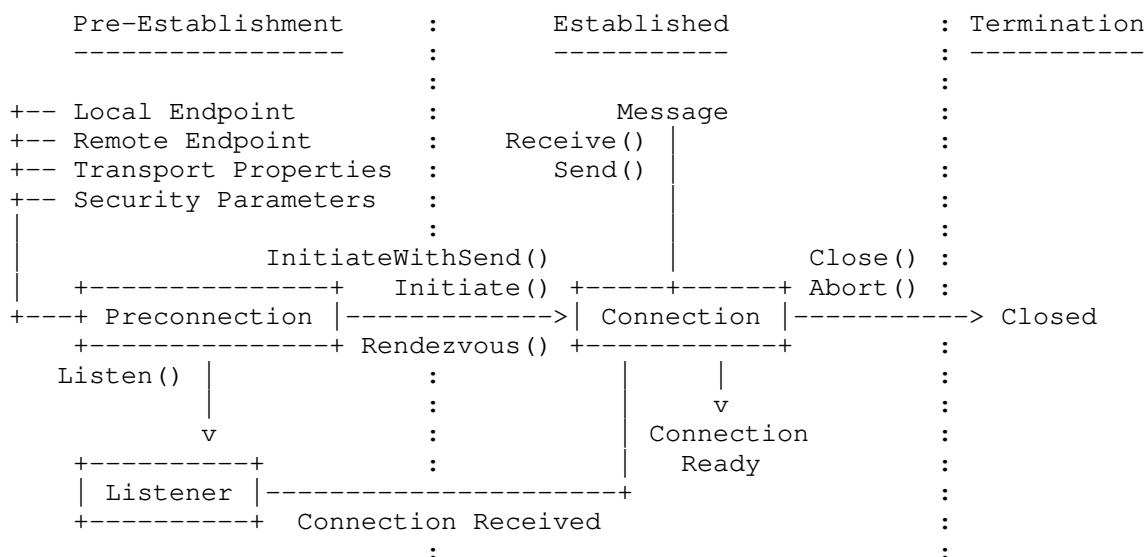


Figure 4: The lifetime of a Connection object

4.1.1. Endpoint Objects

- * **Endpoint:** An Endpoint represents an identifier for one side of a transport connection. Endpoints can be Local Endpoints or Remote Endpoints, and respectively represent an identity that the application uses for the source or destination of a connection. An Endpoint can be specified at various levels of abstraction. An Endpoint at a higher level of abstraction (such as a hostname) can be resolved to more concrete identities (such as IP addresses). An endpoint may also represent a multicast group, in which case it selects a multicast transport for communication.
- * **Remote Endpoint:** The Remote Endpoint represents the application's identifier for a peer that can participate in a transport connection; for example, the combination of a DNS name for the peer and a service name/port.
- * **Local Endpoint:** The Local Endpoint represents the application's identifier for itself that it uses for transport connections; for example, a local IP address and port.

4.1.2. Connections and Related Objects

- * **Preconnection:** A Preconnection object is a representation of a Connection that has not yet been established. It has state that describes parameters of the Connection: the Local Endpoint from

which that Connection will be established, the Remote Endpoint (Section 4.1.3) to which it will connect, and Transport Properties that influence the paths and protocols a Connection will use. A Preconnection can be either fully specified (representing a single possible Connection), or it can be partially specified (representing a family of possible Connections). The Local Endpoint (Section 4.1.3) is required for a Preconnection used to Listen for incoming Connections, but optional if it is used to Initiate a Connection. The Remote Endpoint is required in a Preconnection that used to Initiate a Connection, but is optional if it is used to Listen for incoming Connections. The Local Endpoint and the Remote Endpoint are both required if a peer-to-peer Rendezvous is to occur based on the Preconnection.

- * **Transport Properties:** Transport Properties allow the application to express their requirements, prohibitions, and preferences and configure a Transport Services system. There are three kinds of Transport Properties:
 - **Selection Properties (Section 4.1.3):** Selection Properties can only be specified on a Preconnection.
 - **Connection Properties (Section 4.1.3):** Connection Properties can be specified on a Preconnection and changed on the Connection.
 - **Message Properties (Section 4.1.5):** Message Properties can be specified as defaults on a Preconnection or a Connection, and can also be specified during data transfer to affect specific Messages.
- * **Connection:** A Connection object represents one or more active transport protocol instances that can send and/or receive Messages between Local and Remote Endpoints. It is an abstraction that represents the communication. The Connection object holds state pertaining to the underlying transport protocol instances and any ongoing data transfers. For example, an active Connection can represent a connection-oriented protocol such as TCP, or can represent a fully-specified 5-tuple for a connectionless protocol such as UDP, where the Connection remains an abstraction at the end points. It can also represent a pool of transport protocol instances, e.g., a set of TCP and QUIC connections to equivalent endpoints, or a stream of a multi-streaming transport protocol instance. Connections can be created from a Preconnection or by a Listener.

- * **Listener:** A Listener object accepts incoming transport protocol connections from Remote Endpoints and generates corresponding Connection objects. It is created from a Preconnection object that specifies the type of incoming Connections it will accept.

4.1.3. Pre-Establishment

- * **Selection Properties:** The Selection Properties consist of the properties that an application can set to influence the selection of paths between the Local and Remote Endpoints, to influence the selection of transport protocols, or to configure the behavior of generic transport protocol features. These properties can take the form of requirements, prohibitions, or preferences. Examples of properties that influence path selection include the interface type (such as a Wi-Fi connection, or a Cellular LTE connection), requirements around the largest Message that can be sent, or preferences for throughput and latency. Examples of properties that influence protocol selection and configuration of transport protocol features include reliability, multipath support, and fast open support.
- * **Connection Properties:** The Connection Properties are used to configure protocol-specific options and control per-connection behavior of a Transport Services implementation; for example, a protocol-specific Connection Property can express that if TCP is used, the implementation ought to use the User Timeout Option. Note that the presence of such a property does not require that a specific protocol will be used. In general, these properties do not explicitly determine the selection of paths or protocols, but can be used by an implementation during connection establishment. Connection Properties are specified on a Preconnection prior to Connection establishment, and can be modified on the Connection later. Changes made to Connection Properties after Connection establishment take effect on a best-effort basis.
- * **Security Parameters:** Security Parameters define an application's requirements for authentication and encryption on a Connection. They are used by Transport Security protocols (such as those described in [RFC8922]) to establish secure Connections. Examples of parameters that can be set include local identities, private keys, supported cryptographic algorithms, and requirements for validating trust of remote identities. Security Parameters are primarily associated with a Preconnection object, but properties related to identities can be associated directly with endpoints.

4.1.4. Establishment Actions

- * **Initiate:** The primary action that an application can take to create a Connection to a Remote Endpoint, and prepare any required local or remote state to enable the transmission of Messages. For some protocols, this will initiate a client-to-server style handshake; for other protocols, this will just establish local state (e.g., with connectionless protocols such as UDP). The process of identifying options for connecting, such as resolution of the Remote Endpoint, occurs in response to the Initiate call.
- * **Listen:** Enables a listener to accept incoming Connections. The Listener will then create Connection objects as incoming connections are accepted (Section 4.1.6). Listeners by default register with multiple paths, protocols, and Local Endpoints, unless constrained by Selection Properties and/or the specified Local Endpoint(s). Connections can be accepted on any of the available paths or endpoints.
- * **Rendezvous:** The action of establishing a peer-to-peer connection with a Remote Endpoint. It simultaneously attempts to initiate a connection to a Remote Endpoint while listening for an incoming connection from that endpoint. The process of identifying options for the connection, such as resolution of the Remote Endpoint, occurs in response to the Rendezvous call. As with Listeners, the set of local paths and endpoints is constrained by Selection Properties. If successful, the Rendezvous call returns a Connection object to represent the established peer-to-peer connection. The processes by which connections are initiated during a Rendezvous action will depend on the set of Local and Remote Endpoints configured on the Preconnection. For example, if the Local and Remote Endpoints are TCP host candidates, then a TCP simultaneous open [RFC0793] will be performed. However, if the set of Local Endpoints includes server reflexive candidates, such as those provided by STUN, a Rendezvous action will race candidates in the style of the ICE algorithm [RFC8445] to perform NAT binding discovery and initiate a peer-to-peer connection.

4.1.5. Data Transfer Objects and Actions

- * **Message:** A Message object is a unit of data that can be represented as bytes that can be transferred between two endpoints over a transport connection. The bytes within a Message are assumed to be ordered. If an application does not care about the order in which a peer receives two distinct spans of bytes, those spans of bytes are considered independent Messages.
- * **Message Properties:** Message Properties are used to specify details about Message transmission. They can be specified directly on individual Messages, or can be set on a Preconnection or

Connection as defaults. These properties might only apply to how a Message is sent (such as how the transport will treat prioritization and reliability), but can also include properties that specific protocols encode and communicate to the Remote Endpoint. When receiving Messages, Message Properties can contain information about the received Message, such as metadata generated at the receiver and information signalled by the Remote Endpoint. For example, a Message can be marked with a Message Property indicating that it is the final message on a connection.

- * **Send:** The action to transmit a Message over a Connection to the Remote Endpoint. The interface to Send can accept Message Properties specific to how the Message content is to be sent. The status of the Send operation is delivered back to the sending application in an Event (Section 4.1.6).
- * **Receive:** An action that indicates that the application is ready to asynchronously accept a Message over a Connection from a Remote Endpoint, while the Message content itself will be delivered in an Event (Section 4.1.6). The interface to Receive can include Message Properties specific to the Message that is to be delivered to the application.
- * **Framer:** A Framer is a data translation layer that can be added to a Connection to define how application-layer Messages are transmitted over a transport stack. This is particularly relevant when using a protocol that otherwise presents unstructured streams, such as TCP.

4.1.6. Event Handling

The following categories of events can be delivered to an application:

- * **Connection Ready:** Signals to an application that a given Connection is ready to send and/or receive Messages. If the Connection relies on handshakes to establish state between peers, then it is assumed that these steps have been taken.
- * **Connection Closed:** Signals to an application that a given Connection is no longer usable for sending or receiving Messages. The event delivers a reason or error to the application that describes the nature of the termination.
- * **Connection Received:** Signals to an application that a given Listener has received a Connection.

- * **Message Received:** Delivers received Message content to the application, based on a Receive action. This can include an error if the Receive action cannot be satisfied due to the Connection being closed.
- * **Message Sent:** Notifies the application of the status of its Send action. This might indicate a failure if the Message cannot be sent, or an indication that the Message has been processed by the Transport Services system.
- * **Path Properties Changed:** Notifies the application that a property of the Connection has changed that might influence how and where data is sent and/or received.

4.1.7. Termination Actions

- * **Close:** The action an application takes on a Connection to indicate that it no longer intends to send data, is no longer willing to receive data, and that the protocol should signal this state to the Remote Endpoint if the transport protocol allows this. (Note that this is distinct from the concept of "half-closing" a bidirectional connection, such as when a FIN is sent in one direction of a TCP connection. The end of a stream can also be indicated using Message Properties when sending.)
- * **Abort:** The action the application takes on a Connection to indicate a Close and also indicate that a Transport Services system should not attempt to deliver any outstanding data, and immediately drop the connection. This is intended for immediate, usually abnormal, termination of a connection.

4.1.8. Connection Groups

A Connection Group is a set of Connections that shares properties and caches. A Connection Group represents state for managing Connections within a single application, and does not require end-to-end protocol signaling. For multiplexing transport protocols, only Connections within the same Connection Group are allowed to be multiplexed together.

When the API clones an existing Connection, this adds a new Connection to the Connection Group. A change to one of the Connection Properties on any Connection in the Connection Group automatically changes the Connection Property for all others. All Connections in a Connection Group share the same set of Connection Properties except for the Connection Priority. These Connection Properties are said to be entangled.

For multiplexing transport protocols, only Connections within the same Connection Group are allowed to be multiplexed together. Passive Connections can also be added to a Connection Group, e.g., when a Listener receives a new Connection that is just a new stream of an already active multi-streaming protocol instance.

While Connection Groups are managed by the Transport Services system, an application can define Connection Contexts to control caching boundaries, as discussed in Section 4.2.3.

4.2. Transport Services Implementation

This section defines the key concepts of the Transport Services architecture.

- * Transport Service implementaion: This consists of all objects and protocol instances used internally to a system or library to implement the functionality needed to provide a transport service across a network, as required by the abstract interface.
- * Transport Service system: This consists of the Transport Service implementaion and the Transport Services API.
- * Path: Represents an available set of properties that a local endpoint can use to communicate with a Remote Endpoint, such as routes, addresses, and physical and virtual network interfaces.
- * Protocol Instance: A single instance of one protocol, including any state necessary to establish connectivity or send and receive Messages.
- * Protocol Stack: A set of Protocol Instances (including relevant application, security, transport, or Internet protocols) that are used together to establish connectivity or send and receive Messages. A single stack can be simple (a single transport protocol instance over IP), or it can be complex (multiple application protocol streams going through a single security and transport protocol, over IP; or, a multi-path transport protocol over multiple transport sub-flows).
- * Candidate Path: One path that is available to an application and conforms to the Selection Properties and System Policy, of which there can be several. Candidate Paths are identified during the gathering phase (Section 4.2.1) and can be used during the racing phase (Section 4.2.2).

- * **Candidate Protocol Stack:** One Protocol Stack that can be used by an application for a Connection, which there can be several candidates. Candidate Protocol Stacks are identified during the gathering phase (Section 4.2.1) and are started during the racing phase (Section 4.2.2).
- * **System Policy:** Represents the input from an operating system or other global preferences that can constrain or influence how an implementation will gather candidate paths and Protocol Stacks (Section 4.2.1) and race the candidates during establishment (Section 4.2.2). Specific aspects of the System Policy either apply to all Connections or only certain ones, depending on the runtime context and properties of the Connection.
- * **Cached State:** The state and history that the implementation keeps for each set of associated Endpoints that have been used previously. This can include DNS results, TLS session state, previous success and quality of transport protocols over certain paths, as well as other information.

4.2.1. Candidate Gathering

- * **Candidate Path Selection:** Candidate Path Selection represents the act of choosing one or more paths that are available to use based on the Selection Properties and any available Local and Remote Endpoints provided by the application, as well as the policies and heuristics of a Transport Services implementation.
- * **Candidate Protocol Selection:** Candidate Protocol Selection represents the act of choosing one or more sets of Protocol Stacks that are available to use based on the Transport Properties provided by the application, and the heuristics or policies within the Transport Services implementation.

4.2.2. Candidate Racing

Connection establishment attempts for a set of candidates may be performed simultaneously, synchronously, serially, or using some combination of all of these. We refer to this process as racing, borrowing terminology from Happy Eyeballs [RFC8305].

- * **Protocol Option Racing:** Protocol Option Racing is the act of attempting to establish, or scheduling attempts to establish, multiple Protocol Stacks that differ based on the composition of protocols or the options used for protocols.

- * **Path Racing:** Path Racing is the act of attempting to establish, or scheduling attempts to establish, multiple Protocol Stacks that differ based on a selection from the available Paths. Since different Paths will have distinct configurations for local addresses and DNS servers, attempts across different Paths will perform separate DNS resolution steps, which can lead to further racing of the resolved Remote Endpoints.
- * **Remote Endpoint Racing:** Remote Endpoint Racing is the act of attempting to establish, or scheduling attempts to establish, multiple Protocol Stacks that differ based on the specific representation of the Remote Endpoint, such as a particular IP address that was resolved from a DNS hostname.

4.2.3. Separating Connection Contexts

By default, stored properties of the implementation, such as cached protocol state, cached path state, and heuristics, may be shared (e.g. across multiple connections in an application). This provides efficiency and convenience for the application, since the Transport Services system can automatically optimize behavior.

The Transport Services API can allow applications to explicitly define Connection Contexts that force separation of Cached State and Protocol Stacks. For example, a web browser application could use Connection Contexts with separate caches when implementing different tabs. Possible reasons to isolate Connections using separate Connection Contexts include:

- * Privacy concerns about re-using cached protocol state that can lead to linkability. Sensitive state could include TLS session state [RFC8446] and HTTP cookies [RFC6265]. These concerns could be addressed using Connection Contexts with separate caches, such as for different browser tabs.
- * Privacy concerns about allowing Connections to multiplex together, which can tell a Remote Endpoint that all of the Connections are coming from the same application. Using Connection Contexts avoids the Connections being multiplexed in a HTTP/2 or QUIC stream.

5. IANA Considerations

RFC-EDITOR: Please remove this section before publication.

This document has no actions for IANA.

6. Security and Privacy Considerations

The Transport Services architecture does not recommend use of specific security protocols or algorithms. Its goal is to offer ease of use for existing protocols by providing a generic security-related interface. Each provided interface translates to an existing protocol-specific interface provided by supported security protocols. For example, trust verification callbacks are common parts of TLS APIs; a Transport Services API exposes similar functionality [RFC8922].

As described above in Section 3.3, if a Transport Services implementation races between two different Protocol Stacks, both need to use the same security protocols and options. However, a Transport Services implementation can race different security protocols, e.g., if the application explicitly specifies that it considers them equivalent.

The application controls whether information from previous racing attempts, or other information about past communications that was cached by the Transport Services system is used during establishment. This allows applications to make tradeoffs between efficiency (through racing) and privacy (via information that might leak from the cache toward an on-path observer). Some applications have native concepts (e.g. "incognito mode") that align with this functionality.

Applications need to ensure that they use security APIs appropriately. In cases where applications use an interface to provide sensitive keying material, e.g., access to private keys or copies of pre-shared keys (PSKs), key use needs to be validated. For example, applications ought not to use PSK material created for the Encapsulating Security Protocol (ESP, part of IPsec) [RFC4303] with QUIC, and applications ought not to use private keys intended for server authentication as keys for client authentication.

A Transport Services system must not automatically fall back from secure protocols to insecure protocols, or to weaker versions of secure protocols (see Section 3.3). For example, if an application requests a specific version of TLS, but the desired version of TLS is not available, its connection will fail. Applications are thus responsible for implementing security protocol fallback or version fallback by creating multiple Connections, if so desired. Alternatively, the Transport Services API MAY allow applications to specify that fallback to a specific other version of a protocol is allowed by the Transport Services system.

7. Acknowledgements

This work has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreements No. 644334 (NEAT), No. 688421 (MAMI) and No 815178 (5GENESIS).

This work has been supported by Leibniz Prize project funds of DFG - German Research Foundation: Gottfried Wilhelm Leibniz-Preis 2011 (FKZ FE 570/4-1).

This work has been supported by the UK Engineering and Physical Sciences Research Council under grant EP/R04144X/1.

Thanks to Theresa Enhardt, Max Franke, Mirja Kuehlewind, Jonathan Lennox, and Michael Welzl for the discussions and feedback that helped shape the architecture described here. Particular thanks is also due to Philipp S. Tiesel and Christopher A. Wood, who were both co-authors of this architecture specification as it progressed through the TAPS working group. Thanks as well to Stuart Cheshire, Josh Graessley, David Schinazi, and Eric Kinnear for their implementation and design efforts, including Happy Eyeballs, that heavily influenced this work.

8. References

8.1. Normative References

- [I-D.ietf-taps-interface]
Trammell, B., Welzl, M., Enhardt, T., Fairhurst, G., Kuehlewind, M., Perkins, C., Tiesel, P. S., Wood, C. A., Pauly, T., and K. Rose, "An Abstract Application Layer Interface to Transport Services", Work in Progress, Internet-Draft, draft-ietf-taps-interface-14, 3 January 2022, <<https://datatracker.ietf.org/doc/html/draft-ietf-taps-interface-14>>.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/rfc/rfc2119>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/rfc/rfc8174>>.

8.2. Informative References

- [I-D.ietf-taps-impl]
Brunstrom, A., Pauly, T., Enghardt, T., Grinnemo, K., Jones, T., Tiesel, P. S., Perkins, C., and M. Welzl, "Implementing Interfaces to Transport Services", Work in Progress, Internet-Draft, draft-ietf-taps-impl-10, 12 July 2021, <<https://datatracker.ietf.org/doc/html/draft-ietf-taps-impl-10>>.
- [POSIX] "IEEE Std. 1003.1-2008 Standard for Information Technology -- Portable Operating System Interface (POSIX). Open group Technical Standard: Base Specifications, Issue 7", 2008.
- [RFC0793] Postel, J., "Transmission Control Protocol", STD 7, RFC 793, DOI 10.17487/RFC0793, September 1981, <<https://www.rfc-editor.org/rfc/rfc793>>.
- [RFC4303] Kent, S., "IP Encapsulating Security Payload (ESP)", RFC 4303, DOI 10.17487/RFC4303, December 2005, <<https://www.rfc-editor.org/rfc/rfc4303>>.
- [RFC6265] Barth, A., "HTTP State Management Mechanism", RFC 6265, DOI 10.17487/RFC6265, April 2011, <<https://www.rfc-editor.org/rfc/rfc6265>>.
- [RFC7230] Fielding, R., Ed. and J. Reschke, Ed., "Hypertext Transfer Protocol (HTTP/1.1): Message Syntax and Routing", RFC 7230, DOI 10.17487/RFC7230, June 2014, <<https://www.rfc-editor.org/rfc/rfc7230>>.
- [RFC7540] Belshe, M., Peon, R., and M. Thomson, Ed., "Hypertext Transfer Protocol Version 2 (HTTP/2)", RFC 7540, DOI 10.17487/RFC7540, May 2015, <<https://www.rfc-editor.org/rfc/rfc7540>>.
- [RFC8095] Fairhurst, G., Ed., Trammell, B., Ed., and M. Kuehlewind, Ed., "Services Provided by IETF Transport Protocols and Congestion Control Mechanisms", RFC 8095, DOI 10.17487/RFC8095, March 2017, <<https://www.rfc-editor.org/rfc/rfc8095>>.
- [RFC8305] Schinazi, D. and T. Pauly, "Happy Eyeballs Version 2: Better Connectivity Using Concurrency", RFC 8305, DOI 10.17487/RFC8305, December 2017, <<https://www.rfc-editor.org/rfc/rfc8305>>.

- [RFC8445] Keranen, A., Holmberg, C., and J. Rosenberg, "Interactive Connectivity Establishment (ICE): A Protocol for Network Address Translator (NAT) Traversal", RFC 8445, DOI 10.17487/RFC8445, July 2018, <<https://www.rfc-editor.org/rfc/rfc8445>>.
- [RFC8446] Rescorla, E., "The Transport Layer Security (TLS) Protocol Version 1.3", RFC 8446, DOI 10.17487/RFC8446, August 2018, <<https://www.rfc-editor.org/rfc/rfc8446>>.
- [RFC8922] Enghardt, T., Pauly, T., Perkins, C., Rose, K., and C. Wood, "A Survey of the Interaction between Security Protocols and Transport Services", RFC 8922, DOI 10.17487/RFC8922, October 2020, <<https://www.rfc-editor.org/rfc/rfc8922>>.
- [RFC8923] Welzl, M. and S. Gjessing, "A Minimal Set of Transport Services for End Systems", RFC 8923, DOI 10.17487/RFC8923, October 2020, <<https://www.rfc-editor.org/rfc/rfc8923>>.

Authors' Addresses

Tommy Pauly (editor)
Apple Inc.
One Apple Park Way
Cupertino, California 95014,
United States of America

Email: tpauly@apple.com

Brian Trammell (editor)
Google Switzerland GmbH
Gustav-Gull-Platz 1
CH- 8004 Zurich
Switzerland

Email: ietf@trammell.ch

Anna Brunstrom
Karlstad University
Universitetsgatan 2
651 88 Karlstad
Sweden

Email: anna.brunstrom@kau.se

Godred Fairhurst
University of Aberdeen
Fraser Noble Building
Aberdeen, AB24 3UE

Email: gorry@erg.abdn.ac.uk
URI: <http://www.erg.abdn.ac.uk/>

Colin Perkins
University of Glasgow
School of Computing Science
Glasgow G12 8QQ
United Kingdom

Email: csp@csp Perkins.org

TAPS Working Group
Internet-Draft
Intended status: Informational
Expires: 8 September 2022

A. Brunstrom, Ed.
Karlstad University
T. Pauly, Ed.
Apple Inc.
T. Enghardt
Netflix
P. Tiesel
SAP SE
M. Welzl
University of Oslo
7 March 2022

Implementing Interfaces to Transport Services
draft-ietf-taps-impl-12

Abstract

The Transport Services system enables applications to use transport protocols flexibly for network communication and defines a protocol-independent Transport Services Application Programming Interface (API) that is based on an asynchronous, event-driven interaction pattern. This document serves as a guide to implementation on how to build such a system.

Status of This Memo

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

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 <https://datatracker.ietf.org/drafts/current/>.

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 8 September 2022.

Copyright Notice

Copyright (c) 2022 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 (<https://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 Revised BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Revised BSD License.

Table of Contents

1. Introduction	3
2. Implementing Connection Objects	4
3. Implementing Pre-Establishment	5
3.1. Configuration-time errors	5
3.2. Role of system policy	6
4. Implementing Connection Establishment	7
4.1. Structuring Candidates as a Tree	8
4.1.1. Branch Types	10
4.1.2. Branching Order-of-Operations	12
4.1.3. Sorting Branches	14
4.2. Candidate Gathering	15
4.2.1. Gathering Endpoint Candidates	15
4.3. Candidate Racing	17
4.3.1. Simultaneous	17
4.3.2. Staggered	18
4.3.3. Failover	19
4.4. Completing Establishment	19
4.4.1. Determining Successful Establishment	20
4.5. Establishing multiplexed connections	21
4.6. Handling connectionless protocols	21
4.7. Implementing listeners	21
4.7.1. Implementing listeners for Connected Protocols	22
4.7.2. Implementing listeners for Connectionless Protocols	22
4.7.3. Implementing listeners for Multiplexed Protocols	22
5. Implementing Sending and Receiving Data	23
5.1. Sending Messages	23
5.1.1. Message Properties	23
5.1.2. Send Completion	25
5.1.3. Batching Sends	25
5.2. Receiving Messages	25
5.3. Handling of data for fast-open protocols	26
6. Implementing Message Framers	27
6.1. Defining Message Framers	28
6.2. Sender-side Message Framing	29
6.3. Receiver-side Message Framing	30
7. Implementing Connection Management	31

7.1.	Pooled Connection	31
7.2.	Handling Path Changes	32
8.	Implementing Connection Termination	33
9.	Cached State	34
9.1.	Protocol state caches	34
9.2.	Performance caches	35
10.	Specific Transport Protocol Considerations	36
10.1.	TCP	37
10.2.	MPTCP	39
10.3.	UDP	39
10.4.	UDP-Lite	40
10.5.	UDP Multicast Receive	40
10.6.	SCTP	42
11.	IANA Considerations	45
12.	Security Considerations	45
12.1.	Considerations for Candidate Gathering	45
12.2.	Considerations for Candidate Racing	45
13.	Acknowledgements	46
14.	References	46
14.1.	Normative References	46
14.2.	Informative References	47
	Appendix A. API Mapping Template	49
	Appendix B. Additional Properties	50
	B.1. Properties Affecting Sorting of Branches	50
	Appendix C. Reasons for errors	51
	Appendix D. Existing Implementations	52
	Authors' Addresses	52

1. Introduction

The Transport Services architecture [I-D.ietf-taps-arch] defines a system that allows applications to flexibly use transport networking protocols. The API that such a system exposes to applications is defined as the Transport Services API [I-D.ietf-taps-interface]. This API is designed to be generic across multiple transport protocols and sets of protocols features.

This document serves as a guide to implementation on how to build a system that provides a Transport Services API. It is the job of an implementation of a Transport Services system to turn the requests of an application into decisions on how to establish connections, and how to transfer data over those connections once established. The terminology used in this document is based on the Architecture [I-D.ietf-taps-arch].

2. Implementing Connection Objects

The connection objects that are exposed to applications for Transport Services are:

- * the Preconnection, the bundle of Properties that describes the application constraints on, and preferences for, the transport;
- * the Connection, the basic object that represents a flow of data as Messages in either direction between the Local and Remote Endpoints;
- * and the Listener, a passive waiting object that delivers new Connections.

Preconnection objects should be implemented as bundles of properties that an application can both read and write. A Preconnection object influences a Connection only at one point in time: when the Connection is created. Connection objects represent the interface between the application and the implementation to manage transport state, and conduct data transfer. During the process of establishment (Section 4), the Connection will not be bound to a specific transport protocol instance, since multiple candidate Protocol Stacks might be raced.

Once a Preconnection has been used to create an outbound Connection or a Listener, the implementation should ensure that the copy of the properties held by the Connection or Listener cannot be mutated by the application making changes to the original Preconnection object. This may involve the implementation performing a deep-copy, copying the object with all the objects that it references.

Once the Connection is established, Transport Services implementation maps actions and events to the details of the chosen Protocol Stack. For example, the same Connection object may ultimately represent a single instance of one transport protocol (e.g., a TCP connection, a TLS session over TCP, a UDP flow with fully-specified Local and Remote Endpoints, a DTLS session, a SCTP stream, a QUIC stream, or an HTTP/2 stream). The properties held by a Connection or Listener is independent of other connections that are not part of the same Connection Group.

Connection establishment is only a local operation for a Datagram transport (e.g., UDP(-Lite)), which serves to simplify the local send/receive functions and to filter the traffic for the specified addresses and ports [RFC8085].

Once Initiate has been called, the Selection Properties and Endpoint information are immutable (i.e., an application is not able to later modify Selection Properties on the original Preconnection object). Listener objects are created with a Preconnection, at which point their configuration should be considered immutable by the implementation. The process of listening is described in Section 4.7.

3. Implementing Pre-Establishment

During pre-establishment the application specifies one or more Endpoints to be used for communication as well as protocol preferences and constraints via Selection Properties and, if desired, also Connection Properties. Generally, Connection Properties should be configured as early as possible, because they can serve as input to decisions that are made by the implementation (e.g., the Capacity Profile can guide usage of a protocol offering scavenger-type congestion control).

The implementation stores these properties as a part of the Preconnection object for use during connection establishment. For Selection Properties that are not provided by the application, the implementation must use the default values specified in the Transport Services API ([I-D.ietf-taps-interface]).

3.1. Configuration-time errors

The Transport Services system should have a list of supported protocols available, which each have transport features reflecting the capabilities of the protocol. Once an application specifies its Transport Properties, the transport system matches the required and prohibited properties against the transport features of the available protocols.

In the following cases, failure should be detected during pre-establishment:

- * A request by an application for Protocol Properties that cannot be satisfied by any of the available protocols. For example, if an application requires "Configure Reliability per Message", but no such feature is available in any protocol the host running the transport system on the host running the transport system this should result in an error, e.g., when SCTP is not supported by the operating system.
- * A request by an application for Protocol Properties that are in conflict with each other, i.e., the required and prohibited properties cannot be satisfied by the same protocol. For example,

if an application prohibits "Reliable Data Transfer" but then requires "Configure Reliability per Message", this mismatch should result in an error.

To avoid allocating resources that are not finally needed, it is important that configuration-time errors fail as early as possible.

3.2. Role of system policy

The properties specified during pre-establishment have a close relationship to system policy. The implementation is responsible for combining and reconciling several different sources of preferences when establishing Connections. These include, but are not limited to:

1. Application preferences, i.e., preferences specified during the pre-establishment via Selection Properties.
2. Dynamic system policy, i.e., policy compiled from internally and externally acquired information about available network interfaces, supported transport protocols, and current/previous Connections. Examples of ways to externally retrieve policy-support information are through OS-specific statistics/measurement tools and tools that reside on middleboxes and routers.
3. Default implementation policy, i.e., predefined policy by OS or application.

In general, any protocol or path used for a connection must conform to all three sources of constraints. A violation that occurs at any of the policy layers should cause a protocol or path to be considered ineligible for use. For an example of application preferences leading to constraints, an application may prohibit the use of metered network interfaces for a given Connection to avoid user cost. Similarly, the system policy at a given time may prohibit the use of such a metered network interface from the application's process. Lastly, the implementation itself may default to disallowing certain network interfaces unless explicitly requested by the application and allowed by the system.

It is expected that the database of system policies and the method of looking up these policies will vary across various platforms. An implementation should attempt to look up the relevant policies for the system in a dynamic way to make sure it is reflecting an accurate version of the system policy, since the system's policy regarding the application's traffic may change over time due to user or administrative changes.

4. Implementing Connection Establishment

The process of establishing a network connection begins when an application expresses intent to communicate with a Remote Endpoint by calling Initiate. (At this point, any constraints or requirements the application may have on the connection are available from pre-establishment.) The process can be considered complete once there is at least one Protocol Stack that has completed any required setup to the point that it can transmit and receive the application's data.

Connection establishment is divided into two top-level steps: Candidate Gathering, to identify the paths, protocols, and endpoints to use, and Candidate Racing (see Section 4.2.2 of [I-D.ietf-taps-arch]), in which the necessary protocol handshakes are conducted so that the transport system can select which set to use.

This document structures the candidates for racing as a tree as terminological convention. While a tree structure is not the only way in which racing can be implemented, it does ease the illustration of how racing works.

The most simple example of this process might involve identifying the single IP address to which the implementation wishes to connect, using the system's current default path (i.e., using the default interface), and starting a TCP handshake to establish a stream to the specified IP address. However, each step may also differ depending on the requirements of the connection: if the endpoint is defined as a hostname and port, then there may be multiple resolved addresses that are available; there may also be multiple paths available, (in this case using an interface other than the default system interface); and some protocols may not need any transport handshake to be considered "established" (such as UDP), while other connections may utilize layered protocol handshakes, such as TLS over TCP.

Whenever an implementation has multiple options for connection establishment, it can view the set of all individual connection establishment options as a single, aggregate connection establishment. The aggregate set conceptually includes every valid combination of endpoints, paths, and protocols. As an example, consider an implementation that initiates a TCP connection to a hostname + port endpoint, and has two valid interfaces available (Wi-Fi and LTE). The hostname resolves to a single IPv4 address on the Wi-Fi network, and resolves to the same IPv4 address on the LTE network, as well as a single IPv6 address. The aggregate set of connection establishment options can be viewed as follows:

```
Aggregate [Endpoint: www.example.com:80] [Interface: Any] [Protocol: TCP]
|-> [Endpoint: 192.0.2.1:80] [Interface: Wi-Fi] [Protocol: TCP]
|-> [Endpoint: 192.0.2.1:80] [Interface: LTE] [Protocol: TCP]
|-> [Endpoint: 2001:DB8::1:80] [Interface: LTE] [Protocol: TCP]
```

Any one of these sub-entries on the aggregate connection attempt would satisfy the original application intent. The concern of this section is the algorithm defining which of these options to try, when, and in what order.

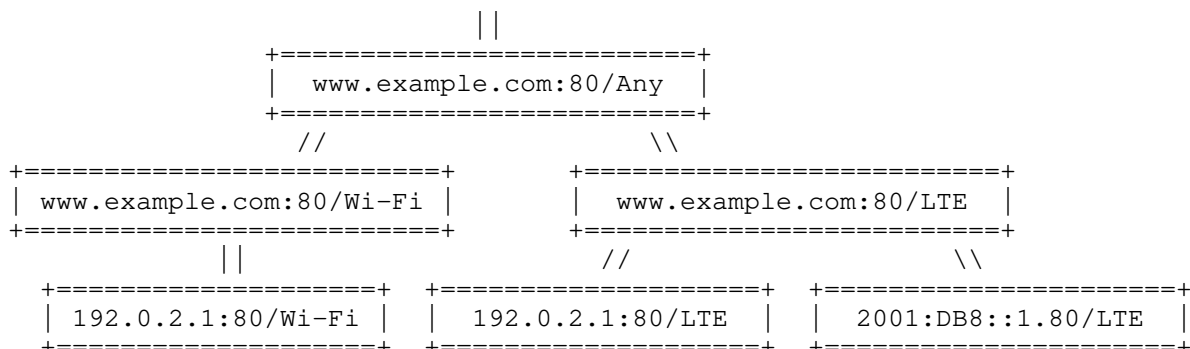
During Candidate Gathering, an implementation first excludes all protocols and paths that match a Prohibit or do not match all Require properties. Then, the implementation will sort branches according to Preferred properties, Avoided properties, and possibly other criteria.

4.1. Structuring Candidates as a Tree

As noted above, the consideration of multiple candidates in a gathering and racing process can be conceptually structured as a tree; this terminological convention is used throughout this document.

Each leaf node of the tree represents a single, coherent connection attempt, with an endpoint, a network path, and a set of protocols that can directly negotiate and send data on the network. Each node in the tree that is not a leaf represents a connection attempt that is either underspecified, or else includes multiple distinct options. For example, when connecting on an IP network, a connection attempt to a hostname and port is underspecified, because the connection attempt requires a resolved IP address as its Remote Endpoint. In this case, the node represented by the connection attempt to the hostname is a parent node, with child nodes for each IP address. Similarly, an implementation that is allowed to connect using multiple interfaces will have a parent node of the tree for the decision between the network paths, with a branch for each interface.

The example aggregate connection attempt above can be drawn as a tree by grouping the addresses resolved on the same interface into branches:



The rest of this section will use a notation scheme to represent this tree. The parent (or trunk) node of the tree will be represented by a single integer, such as "1". Each child of that node will have an integer that identifies it, from 1 to the number of children. That child node will be uniquely identified by concatenating its integer to it's parents identifier with a dot in between, such as "1.1" and "1.2". Each node will be summarized by a tuple of three elements: endpoint, path (labeled here by interface), and protocol. The above example can now be written more succinctly as:

```

1 [www.example.com:80, Any, TCP]
  1.1 [www.example.com:80, Wi-Fi, TCP]
    1.1.1 [192.0.2.1:80, Wi-Fi, TCP]
  1.2 [www.example.com:80, LTE, TCP]
    1.2.1 [192.0.2.1:80, LTE, TCP]
    1.2.2 [2001:DB8::1.80, LTE, TCP]

```

When an implementation views this aggregate set of connection attempts as a single connection establishment, it only will use one of the leaf nodes to transfer data. Thus, when a single leaf node becomes ready to use, then the entire connection attempt is ready to use by the application. Another way to represent this is that every leaf node updates the state of its parent node when it becomes ready, until the trunk node of the tree is ready, which then notifies the application that the connection as a whole is ready to use.

A connection establishment tree may be degenerate, and only have a single leaf node, such as a connection attempt to an IP address over a single interface with a single protocol.

```

1 [192.0.2.1:80, Wi-Fi, TCP]

```

A parent node may also only have one child (or leaf) node, such as a when a hostname resolves to only a single IP address.

```
1 [www.example.com:80, Wi-Fi, TCP]
  1.1 [192.0.2.1:80, Wi-Fi, TCP]
```

4.1.1. Branch Types

There are three types of branching from a parent node into one or more child nodes. Any parent node of the tree must only use one type of branching.

4.1.1.1. Derived Endpoints

If a connection originally targets a single endpoint, there may be multiple endpoints of different types that can be derived from the original. This creates an ordered list of the derived endpoints according to application preference, system policy and expected performance.

DNS hostname-to-address resolution is the most common method of endpoint derivation. When trying to connect to a hostname endpoint on a traditional IP network, the implementation should send DNS queries for both A (IPv4) and AAAA (IPv6) records if both are supported on the local interface. The algorithm for ordering and racing these addresses should follow the recommendations in Happy Eyeballs [RFC8305].

```
1 [www.example.com:80, Wi-Fi, TCP]
  1.1 [2001:DB8::1.80, Wi-Fi, TCP]
  1.2 [192.0.2.1:80, Wi-Fi, TCP]
  1.3 [2001:DB8::2.80, Wi-Fi, TCP]
  1.4 [2001:DB8::3.80, Wi-Fi, TCP]
```

DNS-Based Service Discovery [RFC6763] can also provide an endpoint derivation step. When trying to connect to a named service, the client may discover one or more hostname and port pairs on the local network using multicast DNS [RFC6762]. These hostnames should each be treated as a branch that can be attempted independently from other hostnames. Each of these hostnames might resolve to one or more addresses, which would create multiple layers of branching.

```
1 [term-printer._ipp._tcp.meeting.ietf.org, Wi-Fi, TCP]
  1.1 [term-printer.meeting.ietf.org:631, Wi-Fi, TCP]
    1.1.1 [31.133.160.18.631, Wi-Fi, TCP]
```


Applications can influence which derived endpoints are allowed and preferred via Selection Properties set on the Preconnection. For example, setting a preference for useTemporaryLocalAddress would prefer the use of IPv6 over IPv4, and requiring useTemporaryLocalAddress would eliminate IPv4 options, since IPv4 does not support temporary addresses.

4.1.1.2. Alternate Paths

If a client has multiple network paths available to it, e.g., a mobile client with interfaces for both Wi-Fi and Cellular connectivity, it can attempt a connection over any of the paths. This represents a branch point in the connection establishment. Similar to a derived endpoint, the paths should be ranked based on preference, system policy, and performance. Attempts should be started on one path (e.g., a specific interface), and then successively on other paths (or interfaces) after delays based on expected path round-trip-time or other available metrics.

```
1 [192.0.2.1:80, Any, TCP]
  1.1 [192.0.2.1:80, Wi-Fi, TCP]
  1.2 [192.0.2.1:80, LTE, TCP]
```

This same approach applies to any situation in which the client is aware of multiple links or views of the network. A single interface may be shared by multiple network paths, each with a coherent set of addresses, routes, DNS server, and more. A path may also represent a virtual interface service such as a Virtual Private Network (VPN).

The list of available paths should be constrained by any requirements the application sets, as well as by the system policy.

4.1.1.3. Protocol Options

Differences in possible protocol compositions and options can also provide a branching point in connection establishment. This allows clients to be resilient to situations in which a certain protocol is not functioning on a server or network.

This approach is commonly used for connections with optional proxy server configurations. A single connection might have several options available: an HTTP-based proxy, a SOCKS-based proxy, or no proxy. These options should be ranked and attempted in succession.

- 1 [www.example.com:80, Any, HTTP/TCP]
 - 1.1 [192.0.2.8:80, Any, HTTP/HTTP Proxy/TCP]
 - 1.2 [192.0.2.7:10234, Any, HTTP/SOCKS/TCP]
 - 1.3 [www.example.com:80, Any, HTTP/TCP]
 - 1.3.1 [192.0.2.1:80, Any, HTTP/TCP]

This approach also allows a client to attempt different sets of application and transport protocols that, when available, could provide preferable features. For example, the protocol options could involve QUIC [I-D.ietf-quic-transport] over UDP on one branch, and HTTP/2 [RFC7540] over TLS over TCP on the other:

- 1 [www.example.com:443, Any, Any HTTP]
 - 1.1 [www.example.com:443, Any, QUIC/UDP]
 - 1.1.1 [192.0.2.1:443, Any, QUIC/UDP]
 - 1.2 [www.example.com:443, Any, HTTP2/TLS/TCP]
 - 1.2.1 [192.0.2.1:443, Any, HTTP2/TLS/TCP]

Another example is racing SCTP with TCP:

- 1 [www.example.com:80, Any, Any Stream]
 - 1.1 [www.example.com:80, Any, SCTP]
 - 1.1.1 [192.0.2.1:80, Any, SCTP]
 - 1.2 [www.example.com:80, Any, TCP]
 - 1.2.1 [192.0.2.1:80, Any, TCP]

Implementations that support racing protocols and protocol options should maintain a history of which protocols and protocol options successfully established, on a per-network and per-endpoint basis (see Section 9.2). This information can influence future racing decisions to prioritize or prune branches.

4.1.2. Branching Order-of-Operations

Branch types must occur in a specific order relative to one another to avoid creating leaf nodes with invalid or incompatible settings. In the example above, it would be invalid to branch for derived endpoints (the DNS results for www.example.com) before branching between interface paths, since there are situations when the results will be different across networks due to private names or different supported IP versions. Implementations must be careful to branch in an order that results in usable leaf nodes whenever there are multiple branch types that could be used from a single node.

The order of operations for branching should be:

1. Alternate Paths

2. Protocol Options

3. Derived Endpoints

where a lower number indicates higher precedence and therefore higher placement in the tree. Branching between paths is the first in the list because results across multiple interfaces are likely not related to one another: endpoint resolution may return different results, especially when using locally resolved host and service names, and which protocols are supported and preferred may differ across interfaces. Thus, if multiple paths are attempted, the overall connection can be seen as a race between the available paths or interfaces.

Protocol options are next checked in order. Whether or not a set of protocol, or protocol-specific options, can successfully connect is generally not dependent on which specific IP address is used. Furthermore, the protocol stacks being attempted may influence or altogether change the endpoints being used. Adding a proxy to a connection's branch will change the endpoint to the proxy's IP address or hostname. Choosing an alternate protocol may also modify the ports that should be selected.

Branching for derived endpoints is the final step, and may have multiple layers of derivation or resolution, such as DNS service resolution and DNS hostname resolution.

For example, if the application has indicated both a preference for WiFi over LTE and for a feature only available in SCTP, branches will be first sorted accord to path selection, with WiFi at the top. Then, branches with SCTP will be sorted to the top within their subtree according to the properties influencing protocol selection. However, if the implementation has current cache information that SCTP is not available on the path over WiFi, there is no SCTP node in the WiFi subtree. Here, the path over WiFi will be tried first, and, if connection establishment succeeds, TCP will be used. So the Selection Property of preferring WiFi takes precedence over the Property that led to a preference for SCTP.

```
1. [www.example.com:80, Any, Any Stream]
1.1 [192.0.2.1:80, Wi-Fi, Any Stream]
1.1.1 [192.0.2.1:80, Wi-Fi, TCP]
1.2 [192.0.3.1:80, LTE, Any Stream]
1.2.1 [192.0.3.1:80, LTE, SCTP]
1.2.2 [192.0.3.1:80, LTE, TCP]
```

4.1.3. Sorting Branches

Implementations should sort the branches of the tree of connection options in order of their preference rank, from most preferred to least preferred. Leaf nodes on branches with higher rankings represent connection attempts that will be raced first. Implementations should order the branches to reflect the preferences expressed by the application for its new connection, including Selection Properties, which are specified in [I-D.ietf-taps-interface].

In addition to the properties provided by the application, an implementation may include additional criteria such as cached performance estimates, see Section 9.2, or system policy, see Section 3.2, in the ranking. Two examples of how Selection and Connection Properties may be used to sort branches are provided below:

- * "Interface Instance or Type": If the application specifies an interface type to be preferred or avoided, implementations should accordingly rank the paths. If the application specifies an interface type to be required or prohibited, an implementation is expected to not include the non-conforming paths.
- * "Capacity Profile": An implementation can use the Capacity Profile to prefer paths that match an application's expected traffic pattern. This match will use cached performance estimates, see Section 9.2:
 - Scavenger: Prefer paths with the highest expected available capacity, but minimising impact on other traffic, based on the observed maximum throughput;
 - Low Latency/Interactive: Prefer paths with the lowest expected Round Trip Time, based on observed round trip time estimates;
 - Low Latency/Non-Interactive: Prefer paths with a low expected Round Trip Time, but can tolerate delay variation;
 - Constant-Rate Streaming: Prefer paths that are expected to satisfy the requested Stream Send or Stream Receive Bitrate, based on the observed maximum throughput;
 - Capacity-Seeking: Prefer adapting to paths to determine the highest available capacity, based on the observed maximum throughput.

Implementations process the Properties in the following order: Prohibit, Require, Prefer, Avoid. If Selection Properties contain any prohibited properties, the implementation should first purge branches containing nodes with these properties. For required properties, it should only keep branches that satisfy these requirements. Finally, it should order the branches according to the preferred properties, and finally use any avoided properties as a tiebreaker. When ordering branches, an implementation can give more weight to properties that the application has explicitly set, than to the properties that are default.

The available protocols and paths on a specific system and in a specific context can change; therefore, the result of sorting and the outcome of racing may vary, even when using the same Selection and Connection Properties. However, an implementation ought to provide a consistent outcome to applications, e.g., by preferring protocols and paths that are already used by existing Connections that specified similar Properties.

4.2. Candidate Gathering

The step of gathering candidates involves identifying which paths, protocols, and endpoints may be used for a given Connection. This list is determined by the requirements, prohibitions, and preferences of the application as specified in the Selection Properties.

4.2.1. Gathering Endpoint Candidates

Both Local and Remote Endpoint Candidates must be discovered during connection establishment. To support Interactive Connectivity Establishment (ICE) [RFC8445], or similar protocols that involve out-of-band indirect signalling to exchange candidates with the Remote Endpoint, it is important to query the set of candidate Local Endpoints, and provide the protocol stack with a set of candidate Remote Endpoints, before the Local Endpoint attempts to establish connections.

4.2.1.1. Local Endpoint candidates

The set of possible Local Endpoints is gathered. In the simple case, this merely enumerates the local interfaces and protocols, and allocates ephemeral source ports. For example, a system that has WiFi and Ethernet and supports IPv4 and IPv6 might gather four candidate Local Endpoints (IPv4 on Ethernet, IPv6 on Ethernet, IPv4 on WiFi, and IPv6 on WiFi) that can form the source for a transient.

If NAT traversal is required, the process of gathering Local Endpoints becomes broadly equivalent to the ICE candidate gathering phase (see Section 5.1.1. of [RFC8445]). The endpoint determines its server reflexive Local Endpoints (i.e., the translated address of a Local Endpoint, on the other side of a NAT, e.g via a STUN sever [RFC5389]) and relayed Local Endpoints (e.g., via a TURN server [RFC5766] or other relay), for each interface and network protocol. These are added to the set of candidate Local Endpoints for this connection.

Gathering Local Endpoints is primarily a local operation, although it might involve exchanges with a STUN server to derive server reflexive Local Endpoints, or with a TURN server or other relay to derive relayed Local Endpoints. However, it does not involve communication with the Remote Endpoint.

4.2.1.2. Remote Endpoint Candidates

The Remote Endpoint is typically a name that needs to be resolved into a set of possible addresses that can be used for communication. Resolving the Remote Endpoint is the process of recursively performing such name lookups, until fully resolved, to return the set of candidates for the Remote Endpoint of this connection.

How this resolution is done will depend on the type of the Remote Endpoint, and can also be specific to each Local Endpoint. A common case is when the Remote Endpoint is a DNS name, in which case it is resolved to give a set of IPv4 and IPv6 addresses representing that name. Some types of Remote Endpoint might require more complex resolution. Resolving the Remote Endpoint for a peer-to-peer connection might involve communication with a rendezvous server, which in turn contacts the peer to gain consent to communicate and retrieve its set of candidate Local Endpoints, which are returned and form the candidate remote addresses for contacting that peer.

Resolving the Remote Endpoint is not a local operation. It will involve a directory service, and can require communication with the Remote Endpoint to rendezvous and exchange peer addresses. This can expose some or all of the candidate Local Endpoints to the Remote Endpoint.

4.3. Candidate Racing

The primary goal of the Candidate Racing process is to successfully negotiate a protocol stack to an endpoint over an interface to connect a single leaf node of the tree with as little delay and as few unnecessary connections attempts as possible. Optimizing these two factors improves the user experience, while minimizing network load.

This section covers the dynamic aspect of connection establishment. The tree described above is a useful conceptual and architectural model. However, an implementation is unable to know the full tree before it is formed and many of the possible branches ultimately might not be used.

There are three different approaches to racing the attempts for different nodes of the connection establishment tree:

1. Simultaneous
2. Staggered
3. Failover

Each approach is appropriate in different use-cases and branch types. However, to avoid consuming unnecessary network resources, implementations should not use simultaneous racing as a default approach.

The timing algorithms for racing should remain independent across branches of the tree. Any timers or racing logic is isolated to a given parent node, and is not ordered precisely with regards to other children of other nodes.

4.3.1. Simultaneous

Simultaneous racing is when multiple alternate branches are started without waiting for any one branch to make progress before starting the next alternative. This means the attempts are effectively simultaneous. Simultaneous racing should be avoided by implementations, since it consumes extra network resources and establishes state that might not be used.

4.3.2. Staggered

Staggered racing can be used whenever a single node of the tree has multiple child nodes. Based on the order determined when building the tree, the first child node will be initiated immediately, followed by the next child node after some delay. Once that second child node is initiated, the third child node (if present) will begin after another delay, and so on until all child nodes have been initiated, or one of the child nodes successfully completes its negotiation.

Staggered racing attempts can proceed in parallel. Implementations should not terminate an earlier child connection attempt upon starting a secondary child.

If a child node fails to establish connectivity (as in Section 4.4.1) before the delay time has expired for the next child, the next child should be started immediately.

Staggered racing between IP addresses for a generic Connection should follow the Happy Eyeballs algorithm described in [RFC8305]. [RFC8421] provides guidance for racing when performing Interactive Connectivity Establishment (ICE).

Generally, the delay before starting a given child node ought to be based on the length of time the previously started child node is expected to take before it succeeds or makes progress in connection establishment. Algorithms like Happy Eyeballs choose a delay based on how long the transport connection handshake is expected to take. When performing staggered races in multiple branch types (such as racing between network interfaces, and then racing between IP addresses), a longer delay may be chosen for some branch types. For example, when racing between network interfaces, the delay should also take into account the amount of time it takes to prepare the network interface (such as radio association) and name resolution over that interface, in addition to the delay that would be added for a single transport connection handshake.

Since the staggered delay can be chosen based on dynamic information, such as predicted round-trip time, implementations should define upper and lower bounds for delay times. These bounds are implementation-specific, and may differ based on which branch type is being used.

4.3.3. Failover

If an implementation or application has a strong preference for one branch over another, the branching node may choose to wait until one child has failed before starting the next. Failure of a leaf node is determined by its protocol negotiation failing or timing out; failure of a parent branching node is determined by all of its children failing.

An example in which failover is recommended is a race between a protocol stack that uses a proxy and a protocol stack that bypasses the proxy. Failover is useful in case the proxy is down or misconfigured, but any more aggressive type of racing may end up unnecessarily avoiding a proxy that was preferred by policy.

4.4. Completing Establishment

The process of connection establishment completes when one leaf node of the tree has successfully completed negotiation with the Remote Endpoint, or else all nodes of the tree have failed to connect. The first leaf node to complete its connection is then used by the application to send and receive data.

Successes and failures of a given attempt should be reported up to parent nodes (towards the trunk of the tree). For example, in the following case, if 1.1.1 fails to connect, it reports the failure to 1.1. Since 1.1 has no other child nodes, it also has failed and reports that failure to 1. Because 1.2 has not yet failed, 1 is not considered to have failed. Since 1.2 has not yet started, it is started and the process continues. Similarly, if 1.1.1 successfully connects, then it marks 1.1 as connected, which propagates to the trunk node 1. At this point, the connection as a whole is considered to be successfully connected and ready to process application data.

```
1 [www.example.com:80, Any, TCP]
  1.1 [www.example.com:80, Wi-Fi, TCP]
    1.1.1 [192.0.2.1:80, Wi-Fi, TCP]
    1.2 [www.example.com:80, LTE, TCP]
  ...
```

If a leaf node has successfully completed its connection, all other attempts should be made ineligible for use by the application for the original request. New connection attempts that involve transmitting data on the network ought not to be started after another leaf node has already successfully completed, because the connection as a whole has now been established. An implementation may choose to let certain handshakes and negotiations complete in order to gather metrics to influence future connections. Keeping additional connections is generally not recommended since those attempts were slower to connect and may exhibit less desirable properties.

4.4.1. Determining Successful Establishment

Implementations may select the criteria by which a leaf node is considered to be successfully connected differently on a per-protocol basis. If the only protocol being used is a transport protocol with a clear handshake, like TCP, then the obvious choice is to declare that node "connected" when the last packet of the three-way handshake has been received. If the only protocol being used is an connectionless protocol, like UDP, the implementation may consider the node fully "connected" the moment it determines a route is present, before sending any packets on the network, see further Section 4.6.

For protocol stacks with multiple handshakes, the decision becomes more nuanced. If the protocol stack involves both TLS and TCP, an implementation could determine that a leaf node is connected after the TCP handshake is complete, or it can wait for the TLS handshake to complete as well. The benefit of declaring completion when the TCP handshake finishes, and thus stopping the race for other branches of the tree, is reduced burden on the network and Remote Endpoints from further connection attempts that are likely to be abandoned. On the other hand, by waiting until the TLS handshake is complete, an implementation avoids the scenario in which a TCP handshake completes quickly, but TLS negotiation is either very slow or fails altogether in particular network conditions or to a particular endpoint. To avoid the issue of TLS possibly failing, the implementation should not generate a Ready event for the Connection until TLS is established.

If all of the leaf nodes fail to connect during racing, i.e. none of the configurations that satisfy all requirements given in the Transport Properties actually work over the available paths, then the transport system should notify the application with an `InitiateError` event. An `InitiateError` event should also be generated in case the transport system finds no usable candidates to race.

4.5. Establishing multiplexed connections

Multiplexing several Connections over a single underlying transport connection requires that the Connections to be multiplexed belong to the same Connection Group (as is indicated by the application using the Clone call). When the underlying transport connection supports multi-streaming, the Transport Services System can map each Connection in the Connection Group to a different stream. Thus, when the Connections that are offered to an application by the Transport Services API are multiplexed, the Transport Services implementation can establish a new Connection by simply beginning to use a new stream of an already established transport Connection and there is no need for a connection establishment procedure. This, then, also means that there may not be any "establishment" message (like a TCP SYN), but the application can simply start sending or receiving. Therefore, when the Initiate action of a Transport Services API is called without Messages being handed over, it cannot be guaranteed that the Remote Endpoint will have any way to know about this, and hence a passive endpoint's ConnectionReceived event might not be called until data is received. Instead, calling the ConnectionReceived event could be delayed until the first Message arrives.

4.6. Handling connectionless protocols

While protocols that use an explicit handshake to validate a Connection to a peer can be used for racing multiple establishment attempts in parallel, connectionless protocols such as raw UDP do not offer a way to validate the presence of a peer or the usability of a Connection without application feedback. An implementation should consider such a protocol stack to be established as soon as the Transport Services system has selected a path on which to send data.

However, if a peer is not reachable over the network using the connectionless protocol, or data cannot be exchanged for any other reason, the application may want to attempt using another candidate Protocol Stack. The implementation should maintain the list of other candidate Protocol Stacks that were eligible to use.

4.7. Implementing listeners

When an implementation is asked to Listen, it registers with the system to wait for incoming traffic to the Local Endpoint. If no Local Endpoint is specified, the implementation should use an ephemeral port.

If the Selection Properties do not require a single network interface or path, but allow the use of multiple paths, the Listener object should register for incoming traffic on all of the network interfaces or paths that conform to the Properties. The set of available paths can change over time, so the implementation should monitor network path changes, and change the registration of the Listener across all usable paths as appropriate. When using multiple paths, the Listener is generally expected to use the same port for listening on each.

If the Selection Properties allow multiple protocols to be used for listening, and the implementation supports it, the Listener object should support receiving inbound connections for each eligible protocol on each eligible path.

4.7.1. Implementing listeners for Connected Protocols

Connected protocols such as TCP and TLS-over-TCP have a strong mapping between the Local and Remote Endpoints (four-tuple) and their protocol connection state. These map into Connection objects. Whenever a new inbound handshake is being started, the Listener should generate a new Connection object and pass it to the application.

4.7.2. Implementing listeners for Connectionless Protocols

Connectionless protocols such as UDP and UDP-lite generally do not provide the same mechanisms that connected protocols do to offer Connection objects. Implementations should wait for incoming packets for connectionless protocols on a listening port and should perform four-tuple matching of packets to either existing Connection objects or the creation of new Connection objects. On platforms with facilities to create a "virtual connection" for connectionless protocols implementations should use these mechanisms to minimise the handling of datagrams intended for already created Connection objects.

4.7.3. Implementing listeners for Multiplexed Protocols

Protocols that provide multiplexing of streams into a single four-tuple can listen both for entirely new connections (a new HTTP/2 stream on a new TCP connection, for example) and for new sub-connections (a new HTTP/2 stream on an existing connection). If the abstraction of Connection presented to the application is mapped to the multiplexed stream, then the Listener should deliver new Connection objects in the same way for either case. The implementation should allow the application to introspect the Connection Group marked on the Connections to determine the grouping of the multiplexing.

5. Implementing Sending and Receiving Data

The most basic mapping for sending a Message is an abstraction of datagrams, in which the transport protocol naturally deals in discrete packets. Each Message here corresponds to a single datagram. Generally, these will be short enough that sending and receiving will always use a complete Message.

For protocols that expose byte-streams, the only delineation provided by the protocol is the end of the stream in a given direction. Each Message in this case corresponds to the entire stream of bytes in a direction. These Messages may be quite long, in which case they can be sent in multiple parts.

Protocols that provide the framing (such as length-value protocols, or protocols that use delimiters) may support Message sizes that do not fit within a single datagram. Each Message for framing protocols corresponds to a single frame, which may be sent either as a complete Message in the underlying protocol, or in multiple parts.

5.1. Sending Messages

The effect of the application sending a Message is determined by the top-level protocol in the established Protocol Stack. That is, if the top-level protocol provides an abstraction of framed messages over a connection, the receiving application will be able to obtain multiple Messages on that connection, even if the framing protocol is built on a byte-stream protocol like TCP.

5.1.1. Message Properties

- * **Lifetime:** this should be implemented by removing the Message from the queue of pending Messages after the Lifetime has expired. A queue of pending Messages within the transport system implementation that have yet to be handed to the Protocol Stack can always support this property, but once a Message has been sent into the send buffer of a protocol, only certain protocols may support removing a message. For example, an implementation cannot remove bytes from a TCP send buffer, while it can remove data from a SCTP send buffer using the partial reliability extension [RFC8303]. When there is no standing queue of Messages within the system, and the Protocol Stack does not support the removal of a Message from the stack's send buffer, this property may be ignored.
- * **Priority:** this represents the ability to prioritize a Message over other Messages. This can be implemented by the system re-ordering Messages that have yet to be handed to the Protocol Stack, or by

giving relative priority hints to protocols that support priorities per Message. For example, an implementation of HTTP/2 could choose to send Messages of different Priority on streams of different priority.

- * **Ordered:** when this is false, this disables the requirement of in-order-delivery for protocols that support configurable ordering. When the protocol stack does not support configurable ordering, this property may be ignored.
- * **Safely Replayable:** when this is true, this means that the Message can be used by a transport mechanism that might transfer it multiple times -- e.g., as a result of racing multiple transports or as part of TCP Fast Open. Also, protocols that do not protect against duplicated messages, such as UDP (when used directly, without a protocol layered atop), can only be used with Messages that are Safely Replayable. When a transport system is permitted to replay messages, replay protection could be provided by the application.
- * **Final:** when this is true, this means that the sender will not send any further messages. The Connection need not be closed (in case the Protocol Stack supports half-close operation, like TCP). Any messages sent after a Final message will result in a SendError.
- * **Corruption Protection Length:** when this is set to any value other than Full Coverage, it sets the minimum protection in protocols that allow limiting the checksum length (e.g. UDP-Lite). If the protocol stack does not support checksum length limitation, this property may be ignored.
- * **Reliable Data Transfer (Message):** When true, the property specifies that the Message must be reliably transmitted. When false, and if unreliable transmission is supported by the underlying protocol, then the Message should be unreliably transmitted. If the underlying protocol does not support unreliable transmission, the Message should be reliably transmitted.
- * **Message Capacity Profile Override:** When true, this expresses a wish to override the Generic Connection Property Capacity Profile for this Message. Depending on the value, this can, for example, be implemented by changing the DSCP value of the associated packet (note that the guidelines in Section 6 of [RFC7657] apply; e.g., the DSCP value should not be changed for different packets within a reliable transport protocol session or DCCP connection).

- * No Fragmentation: When set, this property limits the message size to the Maximum Message Size Before Fragmentation or Segmentation (see Section 10.1.7 of [I-D.ietf-taps-interface]). Messages larger than this size generate an error. Setting this avoids transport-layer segmentation or network-layer fragmentation. When used with transports running over IP version 4 the Don't Fragment bit will be set to avoid on-path IP fragmentation ([RFC8304]).

5.1.2. Send Completion

The application should be notified whenever a Message or partial Message has been consumed by the Protocol Stack, or has failed to send. The time at which a Message is considered to have been consumed by the Protocol Stack may vary depending on the protocol. For example, for a basic datagram protocol like UDP, this may correspond to the time when the packet is sent into the interface driver. For a protocol that buffers data in queues, like TCP, this may correspond to when the data has entered the send buffer. The time at which a message failed to send is when Transport Services implementation (including the Protocol Stack) has not successfully sent the entire Message content or partial Message content on any open candidate connection; this can depend on protocol-specific timeouts.

5.1.3. Batching Sends

Since sending a Message may involve a context switch between the application and the Transport Services system, sending patterns that involve multiple small Messages can incur high overhead if each needs to be enqueued separately. To avoid this, the application can indicate a batch of Send actions through the API. When this is used, the implementation can defer the processing of Messages until the batch is complete.

5.2. Receiving Messages

Similar to sending, Receiving a Message is determined by the top-level protocol in the established Protocol Stack. The main difference with Receiving is that the size and boundaries of the Message are not known beforehand. The application can communicate in its Receive action the parameters for the Message, which can help the Transport Services implementation know how much data to deliver and when. For example, if the application only wants to receive a complete Message, the implementation should wait until an entire Message (datagram, stream, or frame) is read before delivering any Message content to the application. This requires the implementation to understand where messages end, either via a supplied deframer or because the top-level protocol in the established Protocol Stack

preserves message boundaries. If the top-level protocol only supports a byte-stream and no framers were supported, the application can control the flow of received data by specifying the minimum number of bytes of Message content it wants to receive at one time.

If a Connection finishes before a requested Receive action can be satisfied, the Transport Services API should deliver any partial Message content outstanding, or if none is available, an indication that there will be no more received Messages.

5.3. Handling of data for fast-open protocols

Several protocols allow sending higher-level protocol or application data during their protocol establishment, such as TCP Fast Open [RFC7413] and TLS 1.3 [RFC8446]. This approach is referred to as sending Zero-RTT (0-RTT) data. This is a desirable feature, but poses challenges to an implementation that uses racing during connection establishment.

The amount of data that can be sent as 0-RTT data varies by protocol and can be queried by the application using the Maximum Message Size Concurrent with Connection Establishment Connection Property. An implementation can set this property according to the protocols that it will race based on the given Selection Properties when the application requests to establish a connection.

If the application has 0-RTT data to send in any protocol handshakes, it needs to provide this data before the handshakes have begun. When racing, this means that the data should be provided before the process of connection establishment has begun. If the application wants to send 0-RTT data, it must indicate this to the implementation by setting the Safely Replayable send parameter to true when sending the data. In general, 0-RTT data may be replayed (for example, if a TCP SYN contains data, and the SYN is retransmitted, the data will be retransmitted as well but may be considered as a new connection instead of a retransmission). Also, when racing connections, different leaf nodes have the opportunity to send the same data independently. If data is truly safely replayable, this should be permissible.

Once the application has provided its 0-RTT data, a Transport Services implementation should keep a copy of this data and provide it to each new leaf node that is started and for which a 0-RTT protocol is being used.

It is also possible that protocol stacks within a particular leaf node use 0-RTT handshakes without any safely replayable application data. For example, TCP Fast Open could use a Client Hello from TLS as its 0-RTT data, shortening the cumulative handshake time.

0-RTT handshakes often rely on previous state, such as TCP Fast Open cookies, previously established TLS tickets, or out-of-band distributed pre-shared keys (PSKs). Implementations should be aware of security concerns around using these tokens across multiple addresses or paths when racing. In the case of TLS, any given ticket or PSK should only be used on one leaf node, since servers will likely reject duplicate tickets in order to prevent replays (see section-8.1 [RFC8446]). If implementations have multiple tickets available from a previous connection, each leaf node attempt can use a different ticket. In effect, each leaf node will send the same early application data, yet encoded (encrypted) differently on the wire.

6. Implementing Message Framers

Message Framers are functions that define simple transformations between application Message data and raw transport protocol data. A Framers can encapsulate or encode outbound Messages, and decapsulate or decode inbound data into Messages.

While many protocols can be represented as Message Framers, for the purposes of the Transport Services API, these are ways for applications or application frameworks to define their own Message parsing to be included within a Connection's Protocol Stack. As an example, TLS is exposed as a protocol natively supported by the Transport Services API, even though it could also serve the purpose of framing data over TCP.

Most Message Framers fall into one of two categories:

- * Header-prefixed record formats, such as a basic Type-Length-Value (TLV) structure
- * Delimiter-separated formats, such as HTTP/1.1.

Common Message Framers can be provided by a Transport Services implementation, but an implementation ought to allow custom Message Framers to be defined by the application or some other piece of software. This section describes one possible API for defining Message Framers as an example.

6.1. Defining Message Framers

A Message Framer is primarily defined by the code that handles events for a framer implementation, specifically how it handles inbound and outbound data parsing. The function that implements custom framing logic will be referred to as the "framer implementation", which may be provided by a Transport Services implementation or the application itself. The Message Framer refers to the object or function within the main Connection implementation that delivers events to the custom framer implementation whenever data is ready to be parsed or framed.

The Transport Services implementation needs to ensure that all of the events and actions taken on a Message Framer are synchronized to ensure consistent behavior. For example, some of the actions defined below (such as `PrependFramer` and `StartPassthrough`) modify how data flows in a protocol stack, and require synchronization with sending and parsing data in the Message Framer.

When a Connection establishment attempt begins, an event can be delivered to notify the framer implementation that a new Connection is being created. Similarly, a stop event can be delivered when a Connection is being torn down. The framer implementation can use the Connection object to look up specific properties of the Connection or the network being used that may influence how to frame Messages.

```
MessageFramer -> Start(Connection)
MessageFramer -> Stop(Connection)
```

When a Message Framer generates a Start event, the framer implementation has the opportunity to start writing some data prior to the Connection delivering its Ready event. This allows the implementation to communicate control data to the Remote Endpoint that can be used to parse Messages.

```
MessageFramer.MakeConnectionReady(Connection)
```

Similarly, when a Message Framer generates a Stop event, the framer implementation has the opportunity to write some final data or clear up its local state before the Closed event is delivered to the Application. The framer implementation can indicate that it has finished with this.

```
MessageFramer.MakeConnectionClosed(Connection)
```

At any time if the implementation encounters a fatal error, it can also cause the Connection to fail and provide an error.

```
MessageFramer.FailConnection(Connection, Error)
```

Should the framer implementation deem the candidate selected during racing unsuitable, it can signal this to the Transport Services API by failing the Connection prior to marking it as ready. If there are no other candidates available, the Connection will fail. Otherwise, the Connection will select a different candidate and the Message Framer will generate a new Start event.

Before an implementation marks a Message Framer as ready, it can also dynamically add a protocol or framer above it in the stack. This allows protocols that need to add TLS conditionally, like STARTTLS [RFC3207], to modify the Protocol Stack based on a handshake result.

```
otherFramer := NewMessageFramer()
MessageFramer.PrependFramer(Connection, otherFramer)
```

A Message Framer might also choose to go into a passthrough mode once an initial exchange or handshake has been completed, such as the STARTTLS case mentioned above. This can also be useful for proxy protocols like SOCKS [RFC1928] or HTTP CONNECT [RFC7230]. In such cases, a Message Framer implementation can intercept sending and receiving of messages at first, but then indicate that no more processing is needed.

```
MessageFramer.StartPassthrough()
```

6.2. Sender-side Message Framing

Message Framers generate an event whenever a Connection sends a new Message.

```
MessageFramer -> NewSentMessage<Connection, MessageData, MessageContext, IsEndOfMessage>
```

Upon receiving this event, a framer implementation is responsible for performing any necessary transformations and sending the resulting data back to the Message Framer, which will in turn send it to the next protocol. Implementations SHOULD ensure that there is a way to pass the original data through without copying to improve performance.

```
MessageFramer.Send(Connection, Data)
```

To provide an example, a simple protocol that adds a length as a header would receive the NewSentMessage event, create a data representation of the length of the Message data, and then send a block of data that is the concatenation of the length header and the original Message data.

6.3. Receiver-side Message Framing

In order to parse a received flow of data into Messages, the Message Framer notifies the framer implementation whenever new data is available to parse.

MessageFramer -> HandleReceivedData<Connection>

Upon receiving this event, the framer implementation can inspect the inbound data. The data is parsed from a particular cursor representing the unprocessed data. The application requests a specific amount of data it needs to have available in order to parse. If the data is not available, the parse fails.

MessageFramer.Parse(Connection, MinimumIncompleteLength, MaximumLength) -> (Data, MessageContext, IsEndOfMessage)

The framer implementation can directly advance the receive cursor once it has parsed data to effectively discard data (for example, discard a header once the content has been parsed).

To deliver a Message to the application, the framer implementation can either directly deliver data that it has allocated, or deliver a range of data directly from the underlying transport and simultaneously advance the receive cursor.

MessageFramer.AdvanceReceiveCursor(Connection, Length)

MessageFramer.DeliverAndAdvanceReceiveCursor(Connection, MessageContext, Length, IsEndOfMessage)

MessageFramer.Deliver(Connection, MessageContext, Data, IsEndOfMessage)

Note that MessageFramer.DeliverAndAdvanceReceiveCursor allows the framer implementation to earmark bytes as part of a Message even before they are received by the transport. This allows the delivery of very large Messages without requiring the implementation to directly inspect all of the bytes.

To provide an example, a simple protocol that parses a length as a header value would receive the HandleReceivedData event, and call Parse with a minimum and maximum set to the length of the header field. Once the parse succeeded, it would call AdvanceReceiveCursor with the length of the header field, and then call DeliverAndAdvanceReceiveCursor with the length of the body that was parsed from the header, marking the new Message as complete.

7. Implementing Connection Management

Once a Connection is established, the Transport Services API allows applications to interact with the Connection by modifying or inspecting Connection Properties. A Connection can also generate events in the form of Soft Errors.

The set of Connection Properties that are supported for setting and getting on a Connection are described in [I-D.ietf-taps-interface]. For any properties that are generic, and thus could apply to all protocols being used by a Connection, the Transport Services implementation should store the properties in storage common to all protocols, and notify all protocol instances in the Protocol Stack whenever the properties have been modified by the application. For protocol-specific properties, such as the User Timeout that applies to TCP, the Transport Services implementation only needs to update the relevant protocol instance.

If an error is encountered in setting a property (for example, if the application tries to set a TCP-specific property on a Connection that is not using TCP), the action should fail gracefully. The application may be informed of the error, but the Connection itself should not be terminated.

The Transport Services API should allow protocol instances in the Protocol Stack to pass up arbitrary generic or protocol-specific errors that can be delivered to the application as Soft Errors. These allow the application to be informed of ICMP errors, and other similar events.

7.1. Pooled Connection

For applications that do not need in-order delivery of Messages, the Transport Services implementation may distribute Messages of a single Connection across several underlying transport connections or multiple streams of multi-streaming connections between endpoints, as long as all of these satisfy the Selection Properties. The Transport Services implementation will then hide this connection management and only expose a single Connection object, which we here call a "Pooled Connection". This is in contrast to Connection Groups, which explicitly expose combined treatment of Connections, giving the application control over multiplexing, for example.

Pooled Connections can be useful when the application using the Transport Services system implements a protocol such as HTTP, which employs request/response pairs and does not require in-order delivery of responses. This enables implementations of Transport Services systems to realize transparent connection coalescing, connection migration, and to perform per-message endpoint and path selection by choosing among multiple underlying connections.

7.2. Handling Path Changes

When a path change occurs, e.g., when the IP address of an interface changes or a new interface becomes available, the Transport Services implementation is responsible for notifying the Protocol Instance of the change. The path change may interrupt connectivity on a path for an active connection or provide an opportunity for a transport that supports multipath or migration to adapt to the new paths. Note that, in the model of the Transport Services API, migration is considered a part of multipath connectivity; it is just a limiting policy on multipath usage. If the multipath Selection Property is set to Disabled, migration is disallowed.

For protocols that do not support multipath or migration, the Protocol Instances should be informed of the path change, but should not be forcibly disconnected if the previously used path becomes unavailable. There are many common user scenarios that can lead to a path becoming temporarily unavailable, and then recovering before the transport protocol reaches a timeout error. These are particularly common using mobile devices. Examples include: an Ethernet cable becoming unplugged and then plugged back in; a device losing a Wi-Fi signal while a user is in an elevator, and reattaching when the user leaves the elevator; and a user losing the radio signal while riding a train through a tunnel. If the device is able to rejoin a network with the same IP address, a stateful transport connection can generally resume. Thus, while it is useful for a Protocol Instance to be aware of a temporary loss of connectivity, the Transport Services implementation should not aggressively close connections in these scenarios.

If the Protocol Stack includes a transport protocol that supports multipath connectivity, the Transport Services implementation should also inform the Protocol Instance of potentially new paths that become permissible based on the multipath Selection Property and the multipath-policy Connection Property choices made by the application. A protocol can then establish new subflows over new paths while an active path is still available or, if migration is supported, also after a break has been detected, and should attempt to tear down subflows over paths that are no longer used. The Connection Property multipath-policy of the Transport Services API allows an application

to indicate when and how different paths should be used. However, detailed handling of these policies is still implementation-specific. For example, if the multipath Selection Property is set to active, the decision about when to create a new path or to announce a new path or set of paths to the Remote Endpoint, e.g., in the form of additional IP addresses, is implementation-specific. If the Protocol Stack includes a transport protocol that does not support multipath, but does support migrating between paths, the update to the set of available paths can trigger the connection to be migrated.

In case of Pooled Connections Section 7.1, the Transport Services implementation may add connections over new paths to the pool if permissible based on the multipath policy and Selection Properties. In case a previously used path becomes unavailable, the transport system may disconnect all connections that require this path, but should not disconnect the pooled connection object exposed to the application. The strategy to do so is implementation-specific, but should be consistent with the behavior of multipath transports.

8. Implementing Connection Termination

With TCP, when an application closes a connection, this means that it has no more data to send (but expects all data that has been handed over to be reliably delivered). However, with TCP only, "close" does not mean that the application will stop receiving data. This is related to TCP's ability to support half-closed connections.

SCTP is an example of a protocol that does not support such half-closed connections. Hence, with SCTP, the meaning of "close" is stricter: an application has no more data to send (but expects all data that has been handed over to be reliably delivered), and will also not receive any more data.

Implementing a protocol independent transport system means that the exposed semantics must be the strictest subset of the semantics of all supported protocols. Hence, as is common with all reliable transport protocols, after a Close action, the application can expect to have its reliability requirements honored regarding the data provided to the Transport Services API, but it cannot expect to be able to read any more data after calling Close.

Abort differs from Close only in that no guarantees are given regarding any data that the application sent to the Transport Services API before calling Abort.

As explained in Section 4.5, when a new stream is multiplexed on an already existing connection of a Transport Protocol Instance, there is no need for a connection establishment procedure. Because the

Connections that are offered by a Transport Services implementation can be implemented as streams that are multiplexed on a transport protocol's connection, it can therefore not be guaranteed an Initiate action from one endpoint provokes a ConnectionReceived event at its peer.

For Close (provoking a Finished event) and Abort (provoking a ConnectionError event), the same logic applies: while it is desirable to be informed when a peer closes or aborts a Connection, whether this is possible depends on the underlying protocol, and no guarantees can be given. With SCTP, the transport system can use the stream reset procedure to cause a Finish event upon a Close action from the peer [NEAT-flow-mapping].

9. Cached State

Beyond a single Connection's lifetime, it is useful for an implementation to keep state and history. This cached state can help improve future Connection establishment due to re-using results and credentials, and favoring paths and protocols that performed well in the past.

Cached state may be associated with different endpoints for the same Connection, depending on the protocol generating the cached content. For example, session tickets for TLS are associated with specific endpoints, and thus should be cached based on a Connection's hostname endpoint (if applicable). However, performance characteristics of a path are more likely tied to the IP address and subnet being used.

9.1. Protocol state caches

Some protocols will have long-term state to be cached in association with endpoints. This state often has some time after which it is expired, so the implementation should allow each protocol to specify an expiration for cached content.

Examples of cached protocol state include:

- * The DNS protocol can cache resolution answers (A and AAAA queries, for example), associated with a Time To Live (TTL) to be used for future hostname resolutions without requiring asking the DNS resolver again.
- * TLS caches session state and tickets based on a hostname, which can be used for resuming sessions with a server.
- * TCP can cache cookies for use in TCP Fast Open.

Cached protocol state is primarily used during Connection establishment for a single Protocol Stack, but may be used to influence an implementation's preference between several candidate Protocol Stacks. For example, if two IP address endpoints are otherwise equally preferred, an implementation may choose to attempt a connection to an address for which it has a TCP Fast Open cookie.

Applications can use the Transport Services API to request that a Connection Group maintain a separate cache for protocol state. Connections in the group will not use cached state from connections outside the group, and connections outside the group will not use state cached from connections inside the group. This may be necessary, for example, if application-layer identifiers rotate and clients wish to avoid linkability via trackable TLS tickets or TFO cookies.

9.2. Performance caches

In addition to protocol state, Protocol Instances should provide data into a performance-oriented cache to help guide future protocol and path selection. Some performance information can be gathered generically across several protocols to allow predictive comparisons between protocols on given paths:

- * Observed Round Trip Time
- * Connection Establishment latency
- * Connection Establishment success rate

These items can be cached on a per-address and per-subnet granularity, and averaged between different values. The information should be cached on a per-network basis, since it is expected that different network attachments will have different performance characteristics. Besides Protocol Instances, other system entities may also provide data into performance-oriented caches. This could for instance be signal strength information reported by radio modems like Wi-Fi and mobile broadband or information about the battery-level of the device. Furthermore, the system may cache the observed maximum throughput on a path as an estimate of the available bandwidth.

An implementation should use this information, when possible, to influence preference between candidate paths, endpoints, and protocol options. Eligible options that historically had significantly better performance than others should be selected first when gathering candidates (see Section 4.2) to ensure better performance for the application.

The reasonable lifetime for cached performance values will vary depending on the nature of the value. Certain information, like the connection establishment success rate to a Remote Endpoint using a given protocol stack, can be stored for a long period of time (hours or longer), since it is expected that the capabilities of the Remote Endpoint are not changing very quickly. On the other hand, the Round Trip Time observed by TCP over a particular network path may vary over a relatively short time interval. For such values, the implementation should remove them from the cache more quickly, or treat older values with less confidence/weight.

[I-D.ietf-tcpm-2140bis] provides guidance about sharing of TCP Control Block information between connections on initialization.

10. Specific Transport Protocol Considerations

Each protocol that is supported by a Transport Services implementation should have a well-defined API mapping. API mappings for a protocol are important for Connections in which a given protocol is the "top" of the Protocol Stack. For example, the mapping of the Send function for TCP applies to Connections in which the application directly sends over TCP.

Each protocol has a notion of Connectedness. Possible values for Connectedness are:

- * Connectionless. Connectionless protocols do not establish explicit state between endpoints, and do not perform a handshake during Connection establishment.
- * Connected. Connected protocols establish state between endpoints, and perform a handshake during Connection establishment. The handshake may be 0-RTT to send data or resume a session, but bidirectional traffic is required to confirm connectedness.
- * Multiplexing Connected. Multiplexing Connected protocols share properties with Connected protocols, but also explicitly support opening multiple application-level flows. This means that they can support cloning new Connection objects without a new explicit handshake.

Protocols also define a notion of Data Unit. Possible values for Data Unit are:

- * Byte-stream. Byte-stream protocols do not define any Message boundaries of their own apart from the end of a stream in each direction.

- * Datagram. Datagram protocols define Message boundaries at the same level of transmission, such that only complete (not partial) Messages are supported.
- * Message. Message protocols support Message boundaries that can be sent and received either as complete or partial Messages. Maximum Message lengths can be defined, and Messages can be partially reliable.

Below, terms in capitals with a dot (e.g., "CONNECT.SCTP") refer to the primitives with the same name in section 4 of [RFC8303]. For further implementation details, the description of these primitives in [RFC8303] points to section 3 of [RFC8303] and section 3 of [RFC8304], which refers back to the relevant specifications for each protocol. This back-tracking method applies to all elements of [RFC8923] (see appendix D of [I-D.ietf-taps-interface]): they are listed in appendix A of [RFC8923] with an implementation hint in the same style, pointing back to section 4 of [RFC8303].

This document defines the API mappings for protocols defined in [RFC8923]. Other protocol mappings can be provided as separate documents, following the mapping template Appendix A.

10.1. TCP

Connectedness: Connected

Data Unit: Byte-stream

API mappings for TCP are as follows:

Connection Object: TCP connections between two hosts map directly to Connection objects.

Initiate: CONNECT.TCP. Calling Initiate on a TCP Connection causes it to reserve a local port, and send a SYN to the Remote Endpoint.

InitiateWithSend: CONNECT.TCP with parameter user message. Early safely replayable data is sent on a TCP Connection in the SYN, as TCP Fast Open data.

Ready: A TCP Connection is ready once the three-way handshake is complete.

InitiateError: Failure of CONNECT.TCP. TCP can throw various errors during connection setup. Specifically, it is important to handle a RST being sent by the peer during the handshake.

ConnectionError: Once established, TCP throws errors whenever the connection is disconnected, such as due to receiving a RST from the peer.

Listen: LISTEN.TCP. Calling Listen for TCP binds a local port and prepares it to receive inbound SYN packets from peers.

ConnectionReceived: TCP Listeners will deliver new connections once they have replied to an inbound SYN with a SYN-ACK.

Clone: Calling Clone on a TCP Connection creates a new Connection with equivalent parameters. These Connections, and Connections generated via later calls to Clone on an Established Connection, form a Connection Group. To realize entanglement for these Connections, with the exception of Connection Priority, changing a Connection Property on one of them must affect the Connection Properties of the others too. No guarantees of honoring the Connection Property Connection Priority are given, and thus it is safe for an implementation of a transport system to ignore this property. When it is reasonable to assume that Connections traverse the same path (e.g., when they share the same encapsulation), support for it can also experimentally be implemented using a congestion control coupling mechanism (see for example [TCP-COUPLING] or [RFC3124]).

Send: SEND.TCP. TCP does not on its own preserve Message boundaries. Calling Send on a TCP connection lays out the bytes on the TCP send stream without any other delineation. Any Message marked as Final will cause TCP to send a FIN once the Message has been completely written, by calling CLOSE.TCP immediately upon successful termination of SEND.TCP. Note that transmitting a Message marked as Final should not cause the Closed event to be delivered to the application, as it will still be possible to receive data until the peer closes or aborts the TCP connection.

Receive: With RECEIVE.TCP, TCP delivers a stream of bytes without any Message delineation. All data delivered in the Received or ReceivedPartial event will be part of a single stream-wide Message that is marked Final (unless a Message Framer is used). EndOfMessage will be delivered when the TCP Connection has received a FIN (CLOSE-EVENT.TCP) from the peer. Note that reception of a FIN should not cause the Closed event to be delivered to the application, as it will still be possible for the application to send data.

Close: Calling Close on a TCP Connection indicates that the

Connection should be gracefully closed (CLOSE.TCP) by sending a FIN to the peer. It will then still be possible to receive data until the peer closes or aborts the TCP connection. The Closed event will be issued upon reception of a FIN.

Abort: Calling Abort on a TCP Connection indicates that the Connection should be immediately closed by sending a RST to the peer (ABORT.TCP).

10.2. MPTCP

Connectedness: Connected

Data Unit: Byte-stream

the Transport Services API mappings for MPTCP are identical to TCP. MPTCP adds support for multipath properties, such as "Multipath Transport" and "Policy for using Multipath Transports".

10.3. UDP

Connectedness: Connectionless

Data Unit: Datagram

API mappings for UDP are as follows:

Connection Object: UDP connections represent a pair of specific IP addresses and ports on two hosts.

Initiate: CONNECT.UDP. Calling Initiate on a UDP Connection causes it to reserve a local port, but does not generate any traffic.

InitiateWithSend: Early data on a UDP Connection does not have any special meaning. The data is sent whenever the Connection is Ready.

Ready: A UDP Connection is ready once the system has reserved a local port and has a path to send to the Remote Endpoint.

InitiateError: UDP Connections can only generate errors on initiation due to port conflicts on the local system.

ConnectionError: Once in use, UDP throws "soft errors" (ERROR.UDP(-Lite)) upon receiving ICMP notifications indicating failures in the network.

Listen: LISTEN.UDP. Calling Listen for UDP binds a local port and

prepares it to receive inbound UDP datagrams from peers.

ConnectionReceived: UDP Listeners will deliver new connections once they have received traffic from a new Remote Endpoint.

Clone: Calling Clone on a UDP Connection creates a new Connection with equivalent parameters. The two Connections are otherwise independent.

Send: SEND.UDP(-Lite). Calling Send on a UDP connection sends the data as the payload of a complete UDP datagram. Marking Messages as Final does not change anything in the datagram's contents. Upon sending a UDP datagram, some relevant fields and flags in the IP header can be controlled: DSCP (SET_DSCP.UDP(-Lite)), DF in IPv4 (SET_DF.UDP(-Lite)) and ECN flag (SET_ECN.UDP(-Lite)).

Receive: RECEIVE.UDP(-Lite). UDP only delivers complete Messages to Received, each of which represents a single datagram received in a UDP packet. Upon receiving a UDP datagram, the ECN flag from the IP header can be obtained (GET_ECN.UDP(-Lite)).

Close: Calling Close on a UDP Connection (ABORT.UDP(-Lite)) releases the local port reservation.

Abort: Calling Abort on a UDP Connection (ABORT.UDP(-Lite)) is identical to calling Close.

10.4. UDP-Lite

Connectedness: Connectionless

Data Unit: Datagram

The Transport Services API mappings for UDP-Lite are identical to UDP. Properties that require checksum coverage are not supported by UDP-Lite, such as "Corruption Protection Length", "Full Checksum Coverage on Sending", "Required Minimum Corruption Protection Coverage for Receiving", and "Full Checksum Coverage on Receiving".

10.5. UDP Multicast Receive

Connectedness: Connectionless

Data Unit: Datagram

API mappings for Receiving Multicast UDP are as follows:

Connection Object: Established UDP Multicast Receive connections

represent a pair of specific IP addresses and ports. The "unidirectional receive" transport property is required, and the Local Endpoint must be configured with a group IP address and a port.

Initiate: Calling `Initiate` on a UDP Multicast Receive Connection causes an immediate `InitiateError`. This is an unsupported operation.

InitiateWithSend: Calling `InitiateWithSend` on a UDP Multicast Receive Connection causes an immediate `InitiateError`. This is an unsupported operation.

Ready: A UDP Multicast Receive Connection is ready once the system has received traffic for the appropriate group and port.

InitiateError: UDP Multicast Receive Connections generate an `InitiateError` if `Initiate` is called.

ConnectionError: Once in use, UDP throws "soft errors" (`ERROR.UDP(-Lite)`) upon receiving ICMP notifications indicating failures in the network.

Listen: `LISTEN.UDP`. Calling `Listen` for UDP Multicast Receive binds a local port, prepares it to receive inbound UDP datagrams from peers, and issues a multicast host join. If a Remote Endpoint with an address is supplied, the join is Source-specific Multicast, and the path selection is based on the route to the Remote Endpoint. If a Remote Endpoint is not supplied, the join is Any-source Multicast, and the path selection is based on the outbound route to the group supplied in the Local Endpoint.

There are cases where it is required to open multiple connections for the same address(es). For example, one Connection might be opened for a multicast group to for a multicast control bus, and another application later opens a separate Connection to the same group to send signals to and/or receive signals from the common bus. In such cases, the Transport Services system needs to explicitly enable re-use of the same set of addresses (equivalent to setting `SO_REUSEADDR` in the socket API).

ConnectionReceived: UDP Multicast Receive Listeners will deliver new connections once they have received traffic from a new Remote Endpoint.

Clone: Calling `Clone` on a UDP Multicast Receive Connection creates a new Connection with equivalent parameters. The two Connections are otherwise independent.

Send: `SEND.UDP(-Lite)`. Calling Send on a UDP Multicast Receive connection causes an immediate `SendError`. This is an unsupported operation.

Receive: `RECEIVE.UDP(-Lite)`. The Receive operation in a UDP Multicast Receive connection only delivers complete Messages to Received, each of which represents a single datagram received in a UDP packet. Upon receiving a UDP datagram, the ECN flag from the IP header can be obtained (`GET_ECN.UDP(-Lite)`).

Close: Calling Close on a UDP Multicast Receive Connection (`ABORT.UDP(-Lite)`) releases the local port reservation and leaves the group.

Abort: Calling Abort on a UDP Multicast Receive Connection (`ABORT.UDP(-Lite)`) is identical to calling Close.

10.6. SCTP

Connectedness: Connected

Data Unit: Message

API mappings for SCTP are as follows:

Connection Object: Connection objects can be mapped to an SCTP association or a stream in an SCTP association. Mapping Connection objects to SCTP streams is called "stream mapping" and has additional requirements as follows. The following explanation assumes a client-server communication model.

Stream mapping requires an association to already be in place between the client and the server, and it requires the server to understand that a new incoming stream should be represented as a new Connection Object by the Transport Services system. A new SCTP stream is created by sending an SCTP message with a new stream id. Thus, to implement stream mapping, the Transport Services API MUST provide a newly created Connection Object to the application upon the reception of such a message. The necessary semantics to implement a Transport Services system Close and Abort primitives are provided by the stream reconfiguration (reset) procedure described in [RFC6525]. This also allows to re-use a stream id after resetting ("closing") the stream. To implement this functionality, SCTP stream reconfiguration [RFC6525] MUST be supported by both the client and the server side.

To avoid head-of-line blocking, stream mapping SHOULD only be implemented when both sides support message interleaving [RFC8260]. This allows a sender to schedule transmissions between multiple streams without risking that transmission of a large message on one stream might block transmissions on other streams for a long time.

To avoid conflicts between stream ids, the following procedure is recommended: the first Connection, for which the SCTP association has been created, MUST always use stream id zero. All additional Connections are assigned to unused stream ids in growing order. To avoid a conflict when both endpoints map new Connections simultaneously, the peer which initiated association MUST use even stream ids whereas the remote side MUST map its Connections to odd stream ids. Both sides maintain a status map of the assigned stream ids. Generally, new streams SHOULD consume the lowest available (even or odd, depending on the side) stream id; this rule is relevant when lower ids become available because Connection objects associated with the streams are closed.

SCTP stream mapping as described here has been implemented in a research prototype; a description of this implementation is given in [NEAT-flow-mapping].

Initiate: If this is the only Connection object that is assigned to the SCTP Association or stream mapping is not used, CONNECT.SCTP is called. Else, unless the Selection Property `activeReadBeforeSend` is Preferred or Required, a new stream is used: if there are enough streams available, Initiate is a local operation that assigns a new stream id to the Connection object. The number of streams is negotiated as a parameter of the prior CONNECT.SCTP call, and it represents a trade-off between local resource usage and the number of Connection objects that can be mapped without requiring a reconfiguration signal. When running out of streams, ADD_STREAM.SCTP must be called.

InitiateWithSend: If this is the only Connection object that is assigned to the SCTP association or stream mapping is not used, CONNECT.SCTP is called with the "user message" parameter. Else, a new stream is used (see Initiate for how to handle running out of streams), and this just sends the first message on a new stream.

Ready: Initiate or InitiateWithSend returns without an error, i.e. SCTP's four-way handshake has completed. If an association with the peer already exists, stream mapping is used and enough streams are available, a Connection Object instantly becomes Ready after calling Initiate or InitiateWithSend.

InitiateError: Failure of CONNECT.SCTP.

ConnectionError: TIMEOUT.SCTP or ABORT-EVENT.SCTP.

Listen: LISTEN.SCTP. If an association with the peer already exists and stream mapping is used, Listen just expects to receive a new message with a new stream id (chosen in accordance with the stream id assignment procedure described above).

ConnectionReceived: LISTEN.SCTP returns without an error (a result of successful CONNECT.SCTP from the peer), or, in case of stream mapping, the first message has arrived on a new stream (in this case, Receive is also invoked).

Clone: Calling Clone on an SCTP association creates a new Connection object and assigns it a new stream id in accordance with the stream id assignment procedure described above. If there are not enough streams available, ADD_STREAM.SCTP must be called.

Priority (Connection): When this value is changed, or a Message with Message Property Priority is sent, and there are multiple Connection objects assigned to the same SCTP association, CONFIGURE_STREAM_SCHEDULER.SCTP is called to adjust the priorities of streams in the SCTP association.

Send: SEND.SCTP. Message Properties such as Lifetime and Ordered map to parameters of this primitive.

Receive: RECEIVE.SCTP. The "partial flag" of RECEIVE.SCTP invokes a ReceivedPartial event.

Close: If this is the only Connection object that is assigned to the SCTP association, CLOSE.SCTP is called, and the Closed event will be delivered to the application upon the ensuing CLOSE-EVENT.SCTP. Else, the Connection object is one out of several Connection objects that are assigned to the same SCTP association, and RESET_STREAM.SCTP must be called, which informs the peer that the stream will no longer be used for mapping and can be used by future Initiate, InitiateWithSend or Listen calls. At the peer, the event RESET_STREAM-EVENT.SCTP will fire, which the peer must answer by issuing RESET_STREAM.SCTP too. The resulting local RESET_STREAM-EVENT.SCTP informs the Transport Services system that the stream id can now be re-used by the next Initiate, InitiateWithSend or Listen calls, and invokes a Closed event towards the application.

Abort: If this is the only Connection object that is assigned to the SCTP association, ABORT.SCTP is called. Else, the Connection object is one out of several Connection objects that are assigned to the same SCTP association, and shutdown proceeds as described under Close.

11. IANA Considerations

RFC-EDITOR: Please remove this section before publication.

This document has no actions for IANA.

12. Security Considerations

[I-D.ietf-taps-arch] outlines general security consideration and requirements for any system that implements the Transport Services architecture. [I-D.ietf-taps-interface] provides further discussion on security and privacy implications of the Transport Services API. This document provides additional guidance on implementation specifics for the Transport Services API and as such the security considerations in both of these documents apply. The next two subsections discuss further considerations that are specific to mechanisms specified in this document.

12.1. Considerations for Candidate Gathering

Implementations should avoid downgrade attacks that allow network interference to cause the implementation to select less secure, or entirely insecure, combinations of paths and protocols.

12.2. Considerations for Candidate Racing

See Section 5.3 for security considerations around racing with 0-RTT data.

An attacker that knows a particular device is racing several options during connection establishment may be able to block packets for the first connection attempt, thus inducing the device to fall back to a secondary attempt. This is a problem if the secondary attempts have worse security properties that enable further attacks. Implementations should ensure that all options have equivalent security properties to avoid incentivizing attacks.

Since results from the network can determine how a connection attempt tree is built, such as when DNS returns a list of resolved endpoints, it is possible for the network to cause an implementation to consume significant on-device resources. Implementations should limit the maximum amount of state allowed for any given node, including the number of child nodes, especially when the state is based on results from the network.

13. Acknowledgements

This work has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No. 644334 (NEAT) and No. 815178 (5GENESIS).

This work has been supported by Leibniz Prize project funds of DFG - German Research Foundation: Gottfried Wilhelm Leibniz-Preis 2011 (FKZ FE 570/4-1).

This work has been supported by the UK Engineering and Physical Sciences Research Council under grant EP/R04144X/1.

This work has been supported by the Research Council of Norway under its "Toppforsk" programme through the "OCARINA" project.

Thanks to Colin Perkins, Tom Jones, Karl-Johan Grinnemo, Gorrry Fairhurst, for their contributions to the design of this specification. Thanks also to Stuart Cheshire, Josh Graessley, David Schinazi, and Eric Kinnear for their implementation and design efforts, including Happy Eyeballs, that heavily influenced this work.

14. References

14.1. Normative References

[I-D.ietf-taps-arch]

Pauly, T., Trammell, B., Brunstrom, A., Fairhurst, G., and C. Perkins, "An Architecture for Transport Services", Work in Progress, Internet-Draft, draft-ietf-taps-arch-12, 3 January 2022, <<https://datatracker.ietf.org/doc/html/draft-ietf-taps-arch-12>>.

[I-D.ietf-taps-interface]

Trammell, B., Welzl, M., Enghardt, T., Fairhurst, G., Kuehlewind, M., Perkins, C., Tiesel, P. S., Wood, C. A., Pauly, T., and K. Rose, "An Abstract Application Layer Interface to Transport Services", Work in Progress, Internet-Draft, draft-ietf-taps-interface-14, 3 January 2022, <<https://datatracker.ietf.org/doc/html/draft-ietf-taps-interface-14>>.

[RFC7413] Cheng, Y., Chu, J., Radhakrishnan, S., and A. Jain, "TCP Fast Open", RFC 7413, DOI 10.17487/RFC7413, December 2014, <<https://www.rfc-editor.org/rfc/rfc7413>>.

- [RFC7540] Belshé, M., Peon, R., and M. Thomson, Ed., "Hypertext Transfer Protocol Version 2 (HTTP/2)", RFC 7540, DOI 10.17487/RFC7540, May 2015, <<https://www.rfc-editor.org/rfc/rfc7540>>.
- [RFC8303] Welzl, M., Tuexen, M., and N. Khademi, "On the Usage of Transport Features Provided by IETF Transport Protocols", RFC 8303, DOI 10.17487/RFC8303, February 2018, <<https://www.rfc-editor.org/rfc/rfc8303>>.
- [RFC8304] Fairhurst, G. and T. Jones, "Transport Features of the User Datagram Protocol (UDP) and Lightweight UDP (UDP-Lite)", RFC 8304, DOI 10.17487/RFC8304, February 2018, <<https://www.rfc-editor.org/rfc/rfc8304>>.
- [RFC8305] Schinazi, D. and T. Pauly, "Happy Eyeballs Version 2: Better Connectivity Using Concurrency", RFC 8305, DOI 10.17487/RFC8305, December 2017, <<https://www.rfc-editor.org/rfc/rfc8305>>.
- [RFC8421] Martinsen, P., Reddy, T., and P. Patil, "Guidelines for Multihomed and IPv4/IPv6 Dual-Stack Interactive Connectivity Establishment (ICE)", BCP 217, RFC 8421, DOI 10.17487/RFC8421, July 2018, <<https://www.rfc-editor.org/rfc/rfc8421>>.
- [RFC8446] Rescorla, E., "The Transport Layer Security (TLS) Protocol Version 1.3", RFC 8446, DOI 10.17487/RFC8446, August 2018, <<https://www.rfc-editor.org/rfc/rfc8446>>.
- [RFC8923] Welzl, M. and S. Gjessing, "A Minimal Set of Transport Services for End Systems", RFC 8923, DOI 10.17487/RFC8923, October 2020, <<https://www.rfc-editor.org/rfc/rfc8923>>.

14.2. Informative References

- [I-D.ietf-quic-transport]
Iyengar, J. and M. Thomson, "QUIC: A UDP-Based Multiplexed and Secure Transport", Work in Progress, Internet-Draft, draft-ietf-quic-transport-34, 14 January 2021, <<https://datatracker.ietf.org/doc/html/draft-ietf-quic-transport-34>>.

- [I-D.ietf-tcpm-2140bis]
Touch, J., Welzl, M., and S. Islam, "TCP Control Block Interdependence", Work in Progress, Internet-Draft, draft-ietf-tcpm-2140bis-11, 12 April 2021, <<https://datatracker.ietf.org/doc/html/draft-ietf-tcpm-2140bis-11>>.
- [NEAT-flow-mapping]
"Transparent Flow Mapping for NEAT", IFIP NETWORKING 2017 Workshop on Future of Internet Transport (FIT 2017) , 2017.
- [RFC1928] Leech, M., Ganis, M., Lee, Y., Kuris, R., Koblas, D., and L. Jones, "SOCKS Protocol Version 5", RFC 1928, DOI 10.17487/RFC1928, March 1996, <<https://www.rfc-editor.org/rfc/rfc1928>>.
- [RFC3124] Balakrishnan, H. and S. Seshan, "The Congestion Manager", RFC 3124, DOI 10.17487/RFC3124, June 2001, <<https://www.rfc-editor.org/rfc/rfc3124>>.
- [RFC3207] Hoffman, P., "SMTP Service Extension for Secure SMTP over Transport Layer Security", RFC 3207, DOI 10.17487/RFC3207, February 2002, <<https://www.rfc-editor.org/rfc/rfc3207>>.
- [RFC5389] Rosenberg, J., Mahy, R., Matthews, P., and D. Wing, "Session Traversal Utilities for NAT (STUN)", RFC 5389, DOI 10.17487/RFC5389, October 2008, <<https://www.rfc-editor.org/rfc/rfc5389>>.
- [RFC5766] Mahy, R., Matthews, P., and J. Rosenberg, "Traversal Using Relays around NAT (TURN): Relay Extensions to Session Traversal Utilities for NAT (STUN)", RFC 5766, DOI 10.17487/RFC5766, April 2010, <<https://www.rfc-editor.org/rfc/rfc5766>>.
- [RFC6525] Stewart, R., Tuexen, M., and P. Lei, "Stream Control Transmission Protocol (SCTP) Stream Reconfiguration", RFC 6525, DOI 10.17487/RFC6525, February 2012, <<https://www.rfc-editor.org/rfc/rfc6525>>.
- [RFC6762] Cheshire, S. and M. Krochmal, "Multicast DNS", RFC 6762, DOI 10.17487/RFC6762, February 2013, <<https://www.rfc-editor.org/rfc/rfc6762>>.
- [RFC6763] Cheshire, S. and M. Krochmal, "DNS-Based Service Discovery", RFC 6763, DOI 10.17487/RFC6763, February 2013, <<https://www.rfc-editor.org/rfc/rfc6763>>.

- [RFC7230] Fielding, R., Ed. and J. Reschke, Ed., "Hypertext Transfer Protocol (HTTP/1.1): Message Syntax and Routing", RFC 7230, DOI 10.17487/RFC7230, June 2014, <<https://www.rfc-editor.org/rfc/rfc7230>>.
- [RFC7657] Black, D., Ed. and P. Jones, "Differentiated Services (Diffserv) and Real-Time Communication", RFC 7657, DOI 10.17487/RFC7657, November 2015, <<https://www.rfc-editor.org/rfc/rfc7657>>.
- [RFC8085] Eggert, L., Fairhurst, G., and G. Shepherd, "UDP Usage Guidelines", BCP 145, RFC 8085, DOI 10.17487/RFC8085, March 2017, <<https://www.rfc-editor.org/rfc/rfc8085>>.
- [RFC8260] Stewart, R., Tuexen, M., Loreto, S., and R. Seggellmann, "Stream Schedulers and User Message Interleaving for the Stream Control Transmission Protocol", RFC 8260, DOI 10.17487/RFC8260, November 2017, <<https://www.rfc-editor.org/rfc/rfc8260>>.
- [RFC8445] Keranen, A., Holmberg, C., and J. Rosenberg, "Interactive Connectivity Establishment (ICE): A Protocol for Network Address Translator (NAT) Traversal", RFC 8445, DOI 10.17487/RFC8445, July 2018, <<https://www.rfc-editor.org/rfc/rfc8445>>.
- [TCP-COUPLING]
"ctrlTCP: Reducing Latency through Coupled, Heterogeneous Multi-Flow TCP Congestion Control", IEEE INFOCOM Global Internet Symposium (GI) workshop (GI 2018) , n.d..

Appendix A. API Mapping Template

Any protocol mapping for the Transport Services API should follow a common template.

Connectedness: (Connectionless/Connected/Multiplexing Connected)

Data Unit: (Byte-stream/Datagram/Message)

Connection Object:

Initiate:

InitiateWithSend:

Ready:

InitiateError:
ConnectionError:
Listen:
ConnectionReceived:
Clone:
Send:
Receive:
Close:
Abort:

Appendix B. Additional Properties

This appendix discusses implementation considerations for additional parameters and properties that could be used to enhance transport protocol and/or path selection, or the transmission of messages given a Protocol Stack that implements them. These are not part of the interface, and may be removed from the final document, but are presented here to support discussion within the TAPS working group as to whether they should be added to a future revision of the base specification.

B.1. Properties Affecting Sorting of Branches

In addition to the Protocol and Path Selection Properties discussed in Section 4.1.3, the following properties under discussion can influence branch sorting:

- * **Bounds on Send or Receive Rate:** If the application indicates a bound on the expected Send or Receive bitrate, an implementation may prefer a path that can likely provide the desired bandwidth, based on cached maximum throughput, see Section 9.2. The application may know the Send or Receive Bitrate from metadata in adaptive HTTP streaming, such as MPEG-DASH.
- * **Cost Preferences:** If the application indicates a preference to avoid expensive paths, and some paths are associated with a monetary cost, an implementation should decrease the ranking of such paths. If the application indicates that it prohibits using expensive paths, paths that are associated with a cost should be purged from the decision tree.

Appendix C. Reasons for errors

The Transport Services API [I-D.ietf-taps-interface] allows for the several generic error types to specify a more detailed reason as to why an error occurred. This appendix lists some of the possible reasons.

- * **InvalidConfiguration:** The transport properties and endpoints provided by the application are either contradictory or incomplete. Examples include the lack of a Remote Endpoint on an active open or using a multicast group address while not requesting a unidirectional receive.
- * **NoCandidates:** The configuration is valid, but none of the available transport protocols can satisfy the transport properties provided by the application.
- * **ResolutionFailed:** The remote or local specifier provided by the application can not be resolved.
- * **EstablishmentFailed:** The Transport Services system was unable to establish a transport-layer connection to the Remote Endpoint specified by the application.
- * **PolicyProhibited:** The system policy prevents the transport system from performing the action requested by the application.
- * **NotCloneable:** The protocol stack is not capable of being cloned.
- * **MessageTooLarge:** The message size is too big for the transport system to handle.
- * **ProtocolFailed:** The underlying protocol stack failed.
- * **InvalidMessageProperties:** The message properties are either contradictory to the transport properties or they can not be satisfied by the transport system.
- * **DeframingFailed:** The data that was received by the underlying protocol stack could not be deframed.
- * **ConnectionAborted:** The connection was aborted by the peer.
- * **Timeout:** Delivery of a message was not possible after a timeout.

Appendix D. Existing Implementations

This appendix gives an overview of existing implementations, at the time of writing, of transport systems that are (to some degree) in line with this document.

* Apple's Network.framework:

- Network.framework is a transport-level API built for C, Objective-C, and Swift. It a connect-by-name API that supports transport security protocols. It provides userspace implementations of TCP, UDP, TLS, DTLS, proxy protocols, and allows extension via custom framers.
- Documentation: <https://developer.apple.com/documentation/network> (<https://developer.apple.com/documentation/network>)

* NEAT and NEATPy:

- NEAT is the output of the European H2020 research project "NEAT"; it is a user-space library for protocol-independent communication on top of TCP, UDP and SCTP, with many more features such as a policy manager.
- Code: <https://github.com/NEAT-project/neat> (<https://github.com/NEAT-project/neat>)
- NEAT project: <https://www.neat-project.org> (<https://www.neat-project.org>)
- NEATPy is a Python shim over NEAT which updates the NEAT API to be in line with version 6 of the Transport Services API draft.
- Code: <https://github.com/theagilepadawan/NEATPy> (<https://github.com/theagilepadawan/NEATPy>)

* PyTAPS:

- A TAPS implementation based on Python asyncio, offering protocol-independent communication to applications on top of TCP, UDP and TLS, with support for multicast.
- Code: <https://github.com/fg-inet/python-asyncio-taps> (<https://github.com/fg-inet/python-asyncio-taps>)

Authors' Addresses

Anna Brunstrom (editor)
Karlstad University
Universitetsgatan 2
651 88 Karlstad
Sweden
Email: anna.brunstrom@kau.se

Tommy Pauly (editor)
Apple Inc.
One Apple Park Way
Cupertino, California 95014,
United States of America
Email: tpauly@apple.com

Theresa Enghardt
Netflix
121 Albright Way
Los Gatos, CA 95032,
United States of America
Email: ietf@tenghardt.net

Philipp S. Tiesel
SAP SE
Konrad-Zuse-Ring 10
14469 Potsdam
Germany
Email: philipp@tiesel.net

Michael Welzl
University of Oslo
PO Box 1080 Blindern
0316 Oslo
Norway
Email: michawe@ifi.uio.no

TAPS Working Group
Internet-Draft
Intended status: Standards Track
Expires: 8 September 2022

B. Trammell, Ed.
Google Switzerland GmbH
M. Welzl, Ed.
University of Oslo
T. Enghardt
Netflix
G. Fairhurst
University of Aberdeen
M. Kuehlewind
Ericsson
C. Perkins
University of Glasgow
P. Tiesel
SAP SE
T. Pauly
Apple Inc.
7 March 2022

An Abstract Application Layer Interface to Transport Services
draft-ietf-taps-interface-15

Abstract

This document describes an abstract application programming interface, API, to the transport layer that enables the selection of transport protocols and network paths dynamically at runtime. This API enables faster deployment of new protocols and protocol features without requiring changes to the applications. The specified API follows the Transport Services architecture by providing asynchronous, atomic transmission of messages. It is intended to replace the BSD sockets API as the common interface to the transport layer, in an environment where endpoints could select from multiple interfaces and potential transport protocols.

Status of This Memo

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

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 <https://datatracker.ietf.org/drafts/current/>.

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 8 September 2022.

Copyright Notice

Copyright (c) 2022 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 (<https://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 Revised BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Revised BSD License.

Table of Contents

1. Introduction	4
1.1. Terminology and Notation	5
1.2. Specification of Requirements	7
2. Overview of the API Design	7
3. API Summary	8
3.1. Usage Examples	9
3.1.1. Server Example	9
3.1.2. Client Example	10
3.1.3. Peer Example	12
4. Transport Properties	13
4.1. Transport Property Names	14
4.2. Transport Property Types	15
5. Scope of the API Definition	15
6. Pre-Establishment Phase	16
6.1. Specifying Endpoints	17
6.1.1. Using Multicast Endpoints	19
6.1.2. Constraining Interfaces for Endpoints	19
6.1.3. Endpoint Aliases	20
6.1.4. Endpoint Examples	20
6.1.5. Multicast Examples	21
6.2. Specifying Transport Properties	23
6.2.1. Reliable Data Transfer (Connection)	26
6.2.2. Preservation of Message Boundaries	27
6.2.3. Configure Per-Message Reliability	27
6.2.4. Preservation of Data Ordering	27

6.2.5.	Use 0-RTT Session Establishment with a Safely Replayable Message	27
6.2.6.	Multistream Connections in Group	28
6.2.7.	Full Checksum Coverage on Sending	28
6.2.8.	Full Checksum Coverage on Receiving	28
6.2.9.	Congestion control	29
6.2.10.	Keep alive	29
6.2.11.	Interface Instance or Type	29
6.2.12.	Provisioning Domain Instance or Type	30
6.2.13.	Use Temporary Local Address	31
6.2.14.	Multipath Transport	32
6.2.15.	Advertisement of Alternative Addresses	33
6.2.16.	Direction of communication	33
6.2.17.	Notification of ICMP soft error message arrival	34
6.2.18.	Initiating side is not the first to write	34
6.3.	Specifying Security Parameters and Callbacks	35
6.3.1.	Specifying Security Parameters on a Pre-Connection . .	35
6.3.2.	Connection Establishment Callbacks	37
7.	Establishing Connections	37
7.1.	Active Open: Initiate	38
7.2.	Passive Open: Listen	39
7.3.	Peer-to-Peer Establishment: Rendezvous	40
7.4.	Connection Groups	42
7.5.	Adding and Removing Endpoints on a Connection	44
8.	Managing Connections	44
8.1.	Generic Connection Properties	46
8.1.1.	Required Minimum Corruption Protection Coverage for Receiving	46
8.1.2.	Connection Priority	47
8.1.3.	Timeout for Aborting Connection	47
8.1.4.	Timeout for keep alive packets	47
8.1.5.	Connection Group Transmission Scheduler	48
8.1.6.	Capacity Profile	48
8.1.7.	Policy for using Multipath Transports	50
8.1.8.	Bounds on Send or Receive Rate	51
8.1.9.	Group Connection Limit	51
8.1.10.	Isolate Session	51
8.1.11.	Read-only Connection Properties	52
8.2.	TCP-specific Properties: User Timeout Option (UTO) . . .	53
8.2.1.	Advertised User Timeout	53
8.2.2.	User Timeout Enabled	53
8.2.3.	Timeout Changeable	54
8.3.	Connection Lifecycle Events	54
8.3.1.	Soft Errors	54
8.3.2.	Path change	54
9.	Data Transfer	54
9.1.	Messages and Framers	55
9.1.1.	Message Contexts	55

9.1.2.	Message Framers	55
9.1.3.	Message Properties	58
9.2.	Sending Data	64
9.2.1.	Basic Sending	64
9.2.2.	Send Events	65
9.2.3.	Partial Sends	66
9.2.4.	Batching Sends	66
9.2.5.	Send on Active Open: InitiateWithSend	67
9.2.6.	Priority and the Transport Services API	67
9.3.	Receiving Data	68
9.3.1.	Enqueueing Receives	68
9.3.2.	Receive Events	69
9.3.3.	Receive Message Properties	71
10.	Connection Termination	73
11.	Connection State and Ordering of Operations and Events	74
12.	IANA Considerations	76
13.	Privacy and Security Considerations	76
14.	Acknowledgements	78
15.	References	78
15.1.	Normative References	78
15.2.	Informative References	79
Appendix A.	Implementation Mapping	83
A.1.	Types	83
A.2.	Events and Errors	84
A.3.	Time Duration	84
Appendix B.	Convenience Functions	84
B.1.	Adding Preference Properties	84
B.2.	Transport Property Profiles	84
B.2.1.	reliable-inorder-stream	84
B.2.2.	reliable-message	85
B.2.3.	unreliable-datagram	85
Appendix C.	Relationship to the Minimal Set of Transport Services for End Systems	86
Authors' Addresses	89

1. Introduction

This document specifies an abstract application programming interface (API) that specifies the interface component of the high-level Transport Services architecture defined in [I-D.ietf-taps-arch]. A Transport Services system supports asynchronous, atomic transmission of messages over transport protocols and network paths dynamically selected at runtime, in environments where an endpoint selects from multiple interfaces and potential transport protocols.

Applications that adopt this API will benefit from a wide set of transport features that can evolve over time. This protocol-independent API ensures that the system providing the API can

optimize its behavior based on the application requirements and network conditions, without requiring changes to the applications. This flexibility enables faster deployment of new features and protocols, and can support applications by offering racing and fallback mechanisms, which otherwise need to be separately implemented in each application.

The Transport Services system derives specific path and protocol selection properties and supported transport features from the analysis provided in [RFC8095], [RFC8923], and [RFC8922]. The Transport Services API enables an implementation to dynamically choose a transport protocol rather than statically binding applications to a protocol at compile time. The Transport Services API also provides applications with a way to override transport selection and instantiate a specific stack, e.g., to support servers wishing to listen to a specific protocol. However, forcing a choice to use a specific transport stack is discouraged for general use, because it can reduce portability.

1.1. Terminology and Notation

The Transport Services API is described in terms of

- * Objects with which an application can interact;
- * Actions the application can perform on these Objects;
- * Events, which an Object can send to an application to be processed asynchronously; and
- * Parameters associated with these Actions and Events.

The following notations, which can be combined, are used in this document:

- * An Action that creates an Object:

Object := Action()

- * An Action that creates an array of Objects:

[]Object := Action()

- * An Action that is performed on an Object:

Object.Action()

- * An Object sends an Event:

Object -> Event<>

- * An Action takes a set of Parameters; an Event contains a set of Parameters. Action and Event parameters whose names are suffixed with a question mark are optional.

Action(param0, param1?, ...) / Event<param0, param1, ...>

Objects that are passed as parameters to Actions use call-by-value behavior. Actions associated with no Object are Actions on the API; they are equivalent to Actions on a per-application global context.

Events are sent to the application or application-supplied code (e.g. framers, see Section 9.1.2) for processing; the details of event processing are platform- and implementation-specific.

We also make use of the following basic types:

- * Boolean: Instances take the value true or false.
- * Integer: Instances take positive or negative integer values.
- * Numeric: Instances take positive or negative real number values.
- * Enumeration: A family of types in which each instance takes one of a fixed, predefined set of values specific to a given enumerated type.
- * Tuple: An ordered grouping of multiple value types, represented as a comma-separated list in parentheses, e.g., (Enumeration, Preference). Instances take a sequence of values each valid for the corresponding value type.
- * Array: Denoted []Type, an instance takes a value for each of zero or more elements in a sequence of the given Type. An array may be of fixed or variable length.
- * Collection: An unordered grouping of one or more values of the same type.

For guidance on how these abstract concepts may be implemented in languages in accordance with native design patterns and language and platform features, see Appendix A.

1.2. Specification of Requirements

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in BCP 14 [RFC2119] [RFC8174] when, and only when, they appear in all capitals, as shown here.

2. Overview of the API Design

The design of the API specified in this document is based on a set of principles, themselves an elaboration on the architectural design principles defined in [I-D.ietf-taps-arch]. The API defined in this document provides:

- * A Transport Services system can offer a variety of transport protocols, independent of the Protocol Stacks that will be used at runtime. All common features of these protocol stacks are made available to the application in a transport-independent way to the degree possible. This enables applications written to a single API to make use of transport protocols in terms of the features they provide.
- * A unified API to datagram and stream-oriented transports, allowing use of a common API for connection establishment and closing.
- * Message-orientation, as opposed to stream-orientation, using application-assisted framing and deframing where the underlying transport does not provide these.
- * Asynchronous Connection establishment, transmission, and reception. This allows concurrent operations during establishment and event-driven application interactions with the transport layer;
- * Selection between alternate network paths, using additional information about the networks over which a connection can operate (e.g. Provisioning Domain (PvD) information [RFC7556]) where available.
- * Explicit support for transport-specific features to be applied, should that particular transport be part of a chosen Protocol Stack.
- * Explicit support for security properties as first-order transport features.

- * Explicit support for configuration of cryptographic identities and transport security parameters persistent across multiple Connections.
- * Explicit support for multistreaming and multipath transport protocols, and the grouping of related Connections into Connection Groups through "cloning" of Connections (see Section 7.4). This function allows applications to take full advantage of new transport protocols supporting these features.

3. API Summary

An application primarily interacts with this API through two Objects: Preconnections and Connections. A Preconnection object (Section 6) represents a set of properties and constraints on the selection and configuration of paths and protocols to establish a Connection with an Endpoint. A Connection object represents an instance of a transport Protocol Stack on which data can be sent to and/or received from a Remote Endpoint (i.e., a logical connection that, depending on the kind of transport, can be bi-directional or unidirectional, and that can use a stream protocol or a datagram protocol). Connections are presented consistently to the application, irrespective of whether the underlying transport is connection-less or connection-oriented. Connections can be created from Preconnections in three ways:

- * by initiating the Preconnection (i.e., actively opening, as in a client; Section 7.1),
- * through listening on the Preconnection (i.e., passively opening, as in a server Section 7.2),
- * or rendezvousing on the Preconnection (i.e., peer to peer establishment; Section 7.3).

Once a Connection is established, data can be sent and received on it in the form of Messages. The API supports the preservation of message boundaries both via explicit Protocol Stack support, and via application support through a Message Framing that finds message boundaries in a stream. Messages are received asynchronously through event handlers registered by the application. Errors and other notifications also happen asynchronously on the Connection. It is not necessary for an application to handle all Events; some Events may have implementation-specific default handlers. The application should not assume that ignoring Events (e.g., Errors) is always safe.

3.1. Usage Examples

The following usage examples illustrate how an application might use the Transport Services API to:

- * Act as a server, by listening for incoming connections, receiving requests, and sending responses, see Section 3.1.1.
- * Act as a client, by connecting to a Remote Endpoint using Initiate, sending requests, and receiving responses, see Section 3.1.2.
- * Act as a peer, by connecting to a Remote Endpoint using Rendezvous while simultaneously waiting for incoming Connections, sending Messages, and receiving Messages, see Section 3.1.3.

The examples in this section presume that a transport protocol is available between the Local and Remote Endpoints that provides Reliable Data Transfer, Preservation of Data Ordering, and Preservation of Message Boundaries. In this case, the application can choose to receive only complete messages.

If none of the available transport protocols provides Preservation of Message Boundaries, but there is a transport protocol that provides a reliable ordered byte stream, an application could receive this byte stream as partial Messages and transform it into application-layer Messages. Alternatively, an application might provide a Message Framing, which can transform a sequence of Messages into a byte stream and vice versa (Section 9.1.2).

3.1.1. Server Example

This is an example of how an application might listen for incoming Connections using the Transport Services API, and receive a request, and send a response.

```
LocalSpecifier := NewLocalEndpoint()
LocalSpecifier.WithInterface("any")
LocalSpecifier.WithService("https")

TransportProperties := NewTransportProperties()
TransportProperties.Require(preserve-msg-boundaries)
// Reliable Data Transfer and Preserve Order are Required by default

SecurityParameters := NewSecurityParameters()
SecurityParameters.Set(identity, myIdentity)
SecurityParameters.Set(key-pair, myPrivateKey, myPublicKey)

// Specifying a Remote Endpoint is optional when using Listen()
Preconnection := NewPreconnection(LocalSpecifier,
                                   TransportProperties,
                                   SecurityParameters)

Listener := Preconnection.Listen()

Listener -> ConnectionReceived<Connection>

// Only receive complete messages in a Conn.Received handler
Connection.Receive()

Connection -> Received<messageDataRequest, messageContext>

//---- Receive event handler begin ----
Connection.Send(messageDataResponse)
Connection.Close()

// Stop listening for incoming Connections
// (this example supports only one Connection)
Listener.Stop()
//---- Receive event handler end ----
```

3.1.2. Client Example

This is an example of how an application might open two Connections to a remote application using the Transport Services API, and send a request as well as receive a response on each of them.

```
RemoteSpecifier := NewRemoteEndpoint()
RemoteSpecifier.WithHostname("example.com")
RemoteSpecifier.WithService("https")

TransportProperties := NewTransportProperties()
TransportProperties.Require(preserve-msg-boundaries)
// Reliable Data Transfer and Preserve Order are Required by default

SecurityParameters := NewSecurityParameters()
TrustCallback := NewCallback({
    // Verify identity of the Remote Endpoint, return the result
})
SecurityParameters.SetTrustVerificationCallback(TrustCallback)

// Specifying a local endpoint is optional when using Initiate()
Preconnection := NewPreconnection(RemoteSpecifier,
                                   TransportProperties,
                                   SecurityParameters)

Connection := Preconnection.Initiate()
Connection2 := Connection.Clone()

Connection -> Ready<>
Connection2 -> Ready<>

//----- Ready event handler for any Connection C begin -----
C.Send(messageDataRequest)

// Only receive complete messages
C.Receive()
//----- Ready event handler for any Connection C end -----

Connection -> Received<messageDataResponse, messageContext>
Connection2 -> Received<messageDataResponse, messageContext>

// Close the Connection in a Receive event handler
Connection.Close()
Connection2.Close()

Preconnections are reusable after being used to initiate a
Connection. Hence, for example, after the Connections were closed,
the following would be correct:

//.. carry out adjustments to the Preconnection, if desire
Connection := Preconnection.Initiate()
```

3.1.3. Peer Example

This is an example of how an application might establish a connection with a peer using `Rendezvous()`, send a `Message`, and receive a `Message`.

```
// Configure local candidates: a port on the Local Endpoint
// and via a STUN server
HostCandidate := NewLocalEndpoint()
HostCandidate.WithPort(9876)

StunCandidate := NewLocalEndpoint()
StunCandidate.WithStunServer(address, port, credentials)

LocalCandidates = [HostCandidate, StunCandidate]

// Configure transport and security properties
TransportProperties := ...
SecurityParameters := ...

Preconnection := NewPreconnection(LocalCandidates,
                                   [], // No remote candidates yet
                                   TransportProperties,
                                   SecurityParameters)

// Resolve the LocalCandidates. The Preconnection.Resolve() call
// resolves both local and remote candidates but, since the remote
// candidates have not yet been specified, the ResolvedRemote list
// returned will be empty and is not used.
ResolvedLocal, ResolvedRemote = Preconnection.Resolve()

// ...Send the ResolvedLocal list to peer via signalling channel
// ...Receive a list of RemoteCandidates from peer via
//     signalling channel

Preconnection.AddRemote(RemoteCandidates)
Preconnection.Rendezvous()

Preconnection -> RendezvousDone<Connection>

//---- RendezvousDone event handler begin ----
Connection.Send(messageDataRequest)
Connection.Receive()
//---- RendezvousDone event handler end ----

Connection -> Received<messageDataResponse, messageContext>

// If new remote endpoint candidates are received from the peer over
```

```
// the signalling channel, for example if using Trickle ICE, then add
// them to the Connection:
Connection.AddRemote(NewRemoteCandidates)

// On a PathChange<> events, resolve the local endpoints to see if a
// new local endpoint has become available and, if so, send to the peer
// as a new candidate and add to the connection:
Connection -> PathChange<>

//---- PathChange event handler begin ----
ResolvedLocal, ResolvedRemote = Preconnection.Resolve()
if ResolvedLocal has changed:
    // ...Send the ResolvedLocal list to peer via signalling channel
    // Add the new local endpoints to the connection:
    Connection.AddLocal(ResolvedLocal)
//---- PathChange event handler end ----

// Close the Connection in a Receive event handler
Connection.Close()
```

4. Transport Properties

Each application using the Transport Services API declares its preferences for how the Transport Services system should operate. This is done by using Transport Properties, as defined in [I-D.ietf-taps-arch], at each stage of the lifetime of a connection.

Transport Properties are divided into Selection, Connection, and Message Properties. Selection Properties (see Section 6.2) can only be set during pre-establishment. They are only used to specify which paths and protocol stacks can be used and are preferred by the application. Although Connection Properties (see Section 8.1) can be set during pre-establishment, they may be changed later. They are used to inform decisions made during establishment and to fine-tune the established connection. Calling Initiate on a Preconnection creates an outbound Connection or a Listener, and the Selection Properties remain readable from the Connection or Listener, but become immutable.

The behavior of the selected protocol stack(s) when sending Messages is controlled by Message Properties (see Section 9.1.3).

Selection Properties can be set on Preconnections, and the effect of Selection Properties can be queried on Connections and Messages. Connection Properties can be set on Connections and Preconnections; when set on Preconnections, they act as an initial default for the resulting Connections. Message Properties can be set on Messages, Connections, and Preconnections; when set on the latter two, they act as an initial default for the Messages sent over those Connections,

Note that configuring Connection Properties and Message Properties on Preconnections is preferred over setting them later. Early specification of Connection Properties allows their use as additional input to the selection process. Protocol Specific Properties, which enable configuration of specialized features of a specific protocol, see Section 3.2 of [I-D.ietf-taps-arch], are not used as an input to the selection process, but only support configuration if the respective protocol has been selected.

4.1. Transport Property Names

Transport Properties are referred to by property names. For the purposes of this document, these names are alphanumeric strings in which words may be separated by hyphens. Specifically, the following characters are allowed: lowercase letters a-z, uppercase letters A-Z, digits 0-9, the hyphen -, and the underscore _. These names serve two purposes:

- * Allowing different components of a Transport Services implementation to pass Transport Properties, e.g., between a language frontend and a policy manager, or as a representation of properties retrieved from a file or other storage.
- * Making the code of different Transport Services implementations look similar. While individual programming languages may preclude strict adherence to the aforementioned naming convention (for instance, by prohibiting the use of hyphens in symbols), users interacting with multiple implementations will still benefit from the consistency resulting from the use of visually similar symbols.

Transport Property Names are hierarchically organized in the form [**<Namespace>.<PropertyName>**].

- * The Namespace component **MUST** be empty for well-known, generic properties, i.e., for properties that are not specific to a protocol and are defined in an RFC.

- * Protocol Specific Properties MUST use the protocol acronym as the Namespace, e.g., tcp for TCP specific Transport Properties. For IETF protocols, property names under these namespaces SHOULD be defined in an RFC.
- * Vendor or implementation specific properties MUST use a string identifying the vendor or implementation as the Namespace.

Namespaces for each of the keywords provided in the IANA protocol numbers registry (see <https://www.iana.org/assignments/protocol-numbers/protocol-numbers.xhtml>) are reserved for Protocol Specific Properties and MUST NOT be used for vendor or implementation-specific properties. Avoid using any of the terms listed as keywords in the protocol numbers registry as any part of a vendor- or implementation-specific property name.

4.2. Transport Property Types

Each Transport Property has a one of the basic types described in Section 1.1.

Most Selection Properties (see Section 6.2) are of the Enumeration type, and use the Preference Enumeration, which takes one of five possible values (Prohibit, Avoid, Ignore, Prefer, or Require) denoting the level of preference for a given property during protocol selection.

5. Scope of the API Definition

This document defines a language- and platform-independent API of a Transport Services system. Given the wide variety of languages and language conventions used to write applications that use the transport layer to connect to other applications over the Internet, this independence makes this API necessarily abstract.

There is no interoperability benefit in tightly defining how the API is presented to application programmers across diverse platforms. However, maintaining the "shape" of the abstract API across different platforms reduces the effort for programmers who learn to use the Transport Services API to then apply their knowledge to another platform.

We therefore make the following recommendations:

- * Actions, Events, and Errors in implementations of the Transport Services API SHOULD use the names given for them in the document, subject to capitalization, punctuation, and other typographic conventions in the language of the implementation, unless the implementation itself uses different names for substantially equivalent objects for networking by convention.
- * Transport Services systems SHOULD implement each Selection Property, Connection Property, and Message Context Property specified in this document. The Transport Services API SHOULD be implemented even when in a specific implementation/platform it will always result in no operation, e.g. there is no action when the API specifies a Property that is not available in a transport protocol implemented on a specific platform. For example, if TCP is the only underlying transport protocol, the Message Property `msgOrdered` can be implemented (trivially, as a no-op) as disabling the requirement for ordering will not have any effect on delivery order for Connections over TCP. Similarly, the `msg-lifetime` Message Property can be implemented but ignored, as the description of this Property states that "it is not guaranteed that a Message will not be sent when its Lifetime has expired".
- * Implementations may use other representations for Transport Property Names, e.g., by providing constants, but should provide a straight-forward mapping between their representation and the property names specified here.

6. Pre-Establishment Phase

The Pre-Establishment phase allows applications to specify properties for the Connections that they are about to make, or to query the API about potential Connections they could make.

A Preconnection Object represents a potential Connection. It is a passive Object (a data structure) that merely maintains the state that describes the properties of a Connection that might exist in the future. This state comprises Local Endpoint and Remote Endpoint Objects that denote the endpoints of the potential Connection (see Section 6.1), the Selection Properties (see Section 6.2), any preconfigured Connection Properties (Section 8.1), and the security parameters (see Section 6.3):

```
Preconnection := NewPreconnection([LocalEndpoint,  
                                   []RemoteEndpoint,  
                                   TransportProperties,  
                                   SecurityParameters])
```

At least one Local Endpoint MUST be specified if the Preconnection is used to Listen() for incoming Connections, but the list of Local Endpoints MAY be empty if the Preconnection is used to Initiate() connections. If no Local Endpoint is specified, the Transport Services system will assign an ephemeral local port to the Connection on the appropriate interface(s). At least one Remote Endpoint MUST be specified if the Preconnection is used to Initiate() Connections, but the list of Remote Endpoints MAY be empty if the Preconnection is used to Listen() for incoming Connections. At least one Local Endpoint and one Remote Endpoint MUST be specified if a peer-to-peer Rendezvous() is to occur based on the Preconnection.

If more than one Local Endpoint is specified on a Preconnection, then all the Local Endpoints on the Preconnection MUST represent the same host. For example, they might correspond to different interfaces on a multi-homed host, or they might correspond to local interfaces and a STUN server that can be resolved to a server reflexive address for a Preconnection used to make a peer-to-peer Rendezvous().

If more than one Remote Endpoint is specified on the Preconnection, then all the Remote Endpoints on the Preconnection SHOULD represent the same service. For example, the Remote Endpoints might represent various network interfaces of a host, or a server reflexive address that can be used to reach a host, or a set of hosts that provide equivalent local balanced service.

In most cases, it is expected that a single Remote Endpoint will be specified by name, and a later call to Initiate() on the Preconnection (see Section 7.1) will internally resolve that name to a list of concrete endpoints. Specifying multiple Remote Endpoints on a Preconnection allows applications to override this for more detailed control.

If Message Framers are used (see Section 9.1.2), they MUST be added to the Preconnection during pre-establishment.

6.1. Specifying Endpoints

The transport services API uses the Local Endpoint and Remote Endpoint Objects to refer to the endpoints of a transport connection. Endpoints can be created as either Remote or Local:

```
RemoteSpecifier := NewRemoteEndpoint()  
LocalSpecifier  := NewLocalEndpoint()
```

A single Endpoint Object represents the identity of a network host. That endpoint can be more or less specific depending on which identifiers are set. For example, an Endpoint that only specifies a hostname may in fact end up corresponding to several different IP addresses on different hosts.

An Endpoint Object can be configured with the following identifiers:

- * Hostname (string):

```
RemoteSpecifier.WithHostname("example.com")
```

- * Port (a 16-bit integer):

```
RemoteSpecifier.WithPort(443)
```

- * Service (an identifier that maps to a port; either a the name of a well-known service, or a DNS SRV service name to be resolved):

```
RemoteSpecifier.WithService("https")
```

- * IP address (IPv4 or IPv6 address):

```
RemoteSpecifier.WithIPv4Address(192.0.2.21)
```

```
RemoteSpecifier.WithIPv6Address(2001:db8:4920:e29d:a420:7461:7073:0a)
```

- * Interface name (string), e.g., to qualify link-local or multicast addresses (see Section 6.1.2 for details):

```
LocalSpecifier.WithInterface("en0")
```

Note that an IPv6 address specified with a scope (e.g. 2001:db8:4920:e29d:a420:7461:7073:0a%en0) is equivalent to WithIPv6Address with an unscoped address and WithInterface together.

An Endpoint cannot have multiple identifiers of a same type set. That is, an endpoint cannot have two IP addresses specified. Two separate IP addresses are represented as two Endpoint Objects. If a Preconnection specifies a Remote Endpoint with a specific IP address set, it will only establish Connections to that IP address. If, on the other hand, the Remote Endpoint specifies a hostname but no addresses, the Connection can perform name resolution and attempt using any address derived from the original hostname of the Remote Endpoint. Note that multiple Remote Endpoints can be added to a Preconnection, as discussed in Section 7.5.

The Transport Services system resolves names internally, when the `Initiate()`, `Listen()`, or `Rendezvous()` method is called to establish a Connection. Privacy considerations for the timing of this resolution are given in Section 13.

The `Resolve()` action on a Preconnection can be used by the application to force early binding when required, for example with some Network Address Translator (NAT) traversal protocols (see Section 7.3).

6.1.1. Using Multicast Endpoints

Specifying a multicast group address on a Local Endpoint will indicate to the Transport Services system that the resulting connection will be used to receive multicast messages. The Remote Endpoint can be used to filter incoming multicast from specific senders. Such a Preconnection will only support calling `Listen()`, not `Initiate()`. Calling `Listen()` will cause the Transport Services system to register for receiving multicast, such as issuing an IGMP join [RFC3376] or using MLD for IPV6 [RFC4604]. Any Connections that are accepted from this Listener are receive-only.

Similarly, specifying a multicast group address on the Remote Endpoint will indicate that the resulting connection will be used to send multicast messages, and that the Preconnection will support `Initiate()` but not `Listen()`. Any Connections created this way are send-only.

A `Rendezvous()` call on Preconnections containing group addresses results in an `EstablishmentError` as described in Section 7.3.

See Section 6.1.5 for more examples.

6.1.2. Constraining Interfaces for Endpoints

Note that this API has multiple ways to constrain and prioritize endpoint candidates based on the network interface:

- * Specifying an interface on a `RemoteEndpoint` qualifies the scope of the remote endpoint, e.g., for link-local addresses.
- * Specifying an interface on a `LocalEndpoint` explicitly binds all candidates derived from this endpoint to use the specified interface.

- * Specifying an interface using the interface Selection Property (Section 6.2.11) or indirectly via the pvd Selection Property (Section 6.2.12) influences the selection among the available candidates.

While specifying an interface on an endpoint restricts the candidates available for connection establishment in the Pre-Establishment Phase, the Selection Properties prioritize and constrain the connection establishment.

6.1.3. Endpoint Aliases

An Endpoint can have an alternative definition when using different protocols. For example, a server that supports both TLS/TCP and QUIC may be accessible on two different port numbers depending on which protocol is used.

To support this, Endpoint Objects can specify "aliases". An Endpoint can have multiple aliases set.

```
RemoteSpecifier.AddAlias(AlternateRemoteSpecifier)
```

In order to scope an alias to a specific transport protocol, an Endpoint can specify a protocol identifier.

```
RemoteSpecifier.WithProtocol(QUIC)
```

The following example shows a case where "example.com" has a server running on port 443, with an alternate port of 8443 for QUIC.

```
RemoteSpecifier := NewRemoteEndpoint()  
RemoteSpecifier.WithHostname("example.com")  
RemoteSpecifier.WithPort(443)  
  
QUICRemoteSpecifier := NewRemoteEndpoint()  
QUICRemoteSpecifier.WithHostname("example.com")  
QUICRemoteSpecifier.WithPort(8443)  
QUICRemoteSpecifier.WithProtocol(QUIC)  
  
RemoteSpecifier.AddAlias(QUICRemoteSpecifier)
```

6.1.4. Endpoint Examples

The following examples of Endpoints show common usage patterns.

Specify a Remote Endpoint using a hostname and service name:

```
RemoteSpecifier := NewRemoteEndpoint()  
RemoteSpecifier.WithHostname("example.com")  
RemoteSpecifier.WithService("https")
```

Specify a Remote Endpoint using an IPv6 address and remote port:

```
RemoteSpecifier := NewRemoteEndpoint()  
RemoteSpecifier.WithIPv6Address(2001:db8:4920:e29d:a420:7461:7073:0a)  
RemoteSpecifier.WithPort(443)
```

Specify a Remote Endpoint using an IPv4 address and remote port:

```
RemoteSpecifier := NewRemoteEndpoint()  
RemoteSpecifier.WithIPv4Address(192.0.2.21)  
RemoteSpecifier.WithPort(443)
```

Specify a Local Endpoint using a local interface name and local port:

```
LocalSpecifier := NewLocalEndpoint()  
LocalSpecifier.WithInterface("en0")  
LocalSpecifier.WithPort(443)
```

As an alternative to specifying an interface name for the Local Endpoint, an application can express more fine-grained preferences using the Interface Instance or Type Selection Property, see Section 6.2.11. However, if the application specifies Selection Properties that are inconsistent with the Local Endpoint, this will result in an Error once the application attempts to open a Connection.

Specify a Local Endpoint using a STUN server:

```
LocalSpecifier := NewLocalEndpoint()  
LocalSpecifier.WithStunServer(address, port, credentials)
```

6.1.5. Multicast Examples

Specify a Local Endpoint using an Any-Source Multicast group to join on a named local interface:

```
LocalSpecifier := NewLocalEndpoint()  
LocalSpecifier.WithIPv4Address(233.252.0.0)  
LocalSpecifier.WithInterface("en0")
```

Source-Specific Multicast requires setting both a Local and Remote Endpoint:


```
LocalSpecifier := NewLocalEndpoint()  
LocalSpecifier.WithIPv4Address(232.1.1.1)  
LocalSpecifier.WithInterface("en0")  
  
RemoteSpecifier := NewRemoteEndpoint()  
RemoteSpecifier.WithIPv4Address(192.0.2.22)
```

One common pattern for multicast is to both send and receive multicast. For such cases, an application can set up both a Listener and a Connection. The Listener is only used to accept Connections that receive inbound multicast. The initiated Connection is only used to send multicast.

```
// Prepare multicast Listener
LocalMulticastSpecifier := NewLocalEndpoint()
LocalMulticastSpecifier.WithIPv4Address(233.252.0.0)
LocalMulticastSpecifier.WithPort(5353)
LocalMulticastSpecifier.WithInterface("en0")

TransportProperties := NewTransportProperties()
TransportProperties.Require(preserve-msg-boundaries)
// Reliable Data Transfer and Preserve Order are Required by default

// Specifying a Remote Endpoint is optional when using Listen()
Preconnection := NewPreconnection(LocalMulticastSpecifier,
                                   TransportProperties,
                                   SecurityParameters)

MulticastListener := Preconnection.Listen()

// Handle inbound messages sent to the multicast group
MulticastListener -> ConnectionReceived<MulticastReceiverConnection>
MulticastReceiverConnection.Receive()
MulticastReceiverConnection -> Received<messageDataRequest, messageContext>

// Prepare Connection to send multicast
LocalSpecifier := NewLocalEndpoint()
LocalSpecifier.WithPort(5353)
LocalSpecifier.WithInterface("en0")
RemoteMulticastSpecifier := NewRemoteEndpoint()
RemoteMulticastSpecifier.WithIPv4Address(233.252.0.0)
RemoteMulticastSpecifier.WithPort(5353)
RemoteMulticastSpecifier.WithInterface("en0")

Preconnection2 := NewPreconnection(LocalSpecifier,
                                   RemoteMulticastSpecifier,
                                   TransportProperties,
                                   SecurityParameters)

// Send outbound messages to the multicast group
MulticastSenderConnection := Preconnection.Initiate()
MulticastSenderConnection.Send(messageData)
```

6.2. Specifying Transport Properties

A Preconnection Object holds properties reflecting the application's requirements and preferences for the transport. These include Selection Properties for selecting protocol stacks and paths, as well as Connection Properties and Message Properties for configuration of the detailed operation of the selected Protocol Stacks on a per-Connection and Message level.

The protocol(s) and path(s) selected as candidates during establishment are determined and configured using these properties. Since there could be paths over which some transport protocols are unable to operate, or remote endpoints that support only specific network addresses or transports, transport protocol selection is necessarily tied to path selection. This may involve choosing between multiple local interfaces that are connected to different access networks.

When additional information (such as Provisioning Domain (PvD) information Path information can include network segment PMTU, set of supported DSCPs, expected usage, cost, etc. The usage of this information by the Transport Services System is generally independent of the specific mechanism/protocol used to receive the information (e.g. zero-conf, DHCP, or IPv6 RA). [RFC7556]) is available about the networks over which an endpoint can operate, this can inform the selection between alternate network paths.

Most Selection Properties are represented as Preferences, which can take one of five values:

Preference	Effect
Require	Select only protocols/paths providing the property, fail otherwise
Prefer	Prefer protocols/paths providing the property, proceed otherwise
Ignore	No preference
Avoid	Prefer protocols/paths not providing the property, proceed otherwise
Prohibit	Select only protocols/paths not providing the property, fail otherwise

Table 1: Selection Property Preference Levels

The implementation **MUST** ensure an outcome that is consistent with all application requirements expressed using **Require** and **Prohibit**. While preferences expressed using **Prefer** and **Avoid** influence protocol and path selection as well, outcomes can vary given the same Selection Properties, because the available protocols and paths can differ across systems and contexts. However, implementations are **RECOMMENDED** to seek to provide a consistent outcome to an application, given the same set of Selection Properties.

Note that application preferences can conflict with each other. For example, if an application indicates a preference for a specific path by specifying an interface, but also a preference for a protocol, a situation might occur in which the preferred protocol is not available on the preferred path. In such cases, applications can expect properties that determine path selection to be prioritized over properties that determine protocol selection. The transport system **SHOULD** determine the preferred path first, regardless of protocol preferences. This ordering is chosen to provide consistency across implementations, based on the fact that it is more common for the use of a given network path to determine cost to the user (i.e., an interface type preference might be based on a user's preference to avoid being charged more for a cellular data plan).

Selection and Connection Properties, as well as defaults for Message Properties, can be added to a Preconnection to configure the selection process and to further configure the eventually selected protocol stack(s). They are collected into a TransportProperties object to be passed into a Preconnection object:

```
TransportProperties := NewTransportProperties()
```

Individual properties are then set on the TransportProperties Object. Setting a Transport Property to a value overrides the previous value of this Transport Property.

```
TransportProperties.Set(property, value)
```

To aid readability, implementations **MAY** provide additional convenience functions to simplify use of Selection Properties: see Appendix B.1 for examples. In addition, implementations **MAY** provide a mechanism to create TransportProperties objects that are preconfigured for common use cases as outlined in Appendix B.2.

Transport Properties for an established connection can be queried via the Connection object, as outlined in Section 8.

A Connection gets its Transport Properties either by being explicitly configured via a Preconnection, by configuration after establishment, or by inheriting them from an antecedent via cloning; see Section 7.4 for more.

Section 8.1 provides a list of Connection Properties, while Selection Properties are listed in the subsections below. Many properties are only considered during establishment, and can not be changed after a Connection is established; however, they can still be queried. The return type of a queried Selection Property is Boolean, where true means that the Selection Property has been applied and false means that the Selection Property has not been applied. Note that true does not mean that a request has been honored. For example, if Congestion control was requested with preference level Prefer, but congestion control could not be supported, querying the congestionControl property yields the value false. If the preference level Avoid was used for Congestion control, and, as requested, the Connection is not congestion controlled, querying the congestionControl property also yields the value false.

An implementation of the Transport Services API must provide sensible defaults for Selection Properties. The default values for each property below represent a configuration that can be implemented over TCP. If these default values are used and TCP is not supported by a Transport Services system, then an application using the default set of Properties might not succeed in establishing a connection. Using the same default values for independent Transport Services implementations can be beneficial when applications are ported between different implementations/platforms, even if this default could lead to a connection failure when TCP is not available. If default values other than those suggested below are used, it is RECOMMENDED to clearly document any differences.

6.2.1. Reliable Data Transfer (Connection)

Name: reliability

Type: Preference

Default: Require

This property specifies whether the application needs to use a transport protocol that ensures that all data is received at the Remote Endpoint without corruption. When reliable data transfer is enabled, this also entails being notified when a Connection is closed or aborted.

6.2.2. Preservation of Message Boundaries

Name: preserveMsgBoundaries

Type: Preference

Default: Ignore

This property specifies whether the application needs or prefers to use a transport protocol that preserves message boundaries.

6.2.3. Configure Per-Message Reliability

Name: perMsgReliability

Type: Preference

Default: Ignore

This property specifies whether an application considers it useful to specify different reliability requirements for individual Messages in a Connection.

6.2.4. Preservation of Data Ordering

Name: preserveOrder

Type: Preference

Default: Require

This property specifies whether the application wishes to use a transport protocol that can ensure that data is received by the application on the other end in the same order as it was sent.

6.2.5. Use 0-RTT Session Establishment with a Safely Replayable Message

Name: zeroRttMsg

Type: Preference

Default: Ignore

This property specifies whether an application would like to supply a Message to the transport protocol before Connection establishment that will then be reliably transferred to the other side before or during Connection establishment. This Message can potentially be received multiple times (i.e., multiple copies of the message data may be passed to the Remote Endpoint). See also Section 9.1.3.4.

6.2.6. Multistream Connections in Group

Name: multistreaming

Type: Preference

Default: Prefer

This property specifies that the application would prefer multiple Connections within a Connection Group to be provided by streams of a single underlying transport connection where possible.

6.2.7. Full Checksum Coverage on Sending

Name: fullChecksumSend

Type: Preference

Default: Require

This property specifies the application's need for protection against corruption for all data transmitted on this Connection. Disabling this property could enable later control of the sender checksum coverage (see Section 9.1.3.6).

6.2.8. Full Checksum Coverage on Receiving

Name: fullChecksumRecv

Type: Preference

Default: Require

This property specifies the application's need for protection against corruption for all data received on this Connection. Disabling this property could enable later control of the required minimum receiver checksum coverage (see Section 8.1.1).

6.2.9. Congestion control

Name: congestionControl

Type: Preference

Default: Require

This property specifies whether the application would like the Connection to be congestion controlled or not. Note that if a Connection is not congestion controlled, an application using such a Connection SHOULD itself perform congestion control in accordance with [RFC2914] or use a circuit breaker in accordance with [RFC8084], whichever is appropriate. Also note that reliability is usually combined with congestion control in protocol implementations, rendering "reliable but not congestion controlled" a request that is unlikely to succeed. If the Connection is congestion controlled, performing additional congestion control in the application can have negative performance implications.

6.2.10. Keep alive

Name: keepAlive

Type: Preference

Default: Ignore

This property specifies whether the application would like the Connection to send keep-alive packets or not. Note that if a Connection determines that keep-alive packets are being sent, the application should itself avoid generating additional keep alive messages. Note that when supported, the system will use the default period for generation of the keep alive-packets. (See also Section 8.1.4).

6.2.11. Interface Instance or Type

Name: interface

Type: Collection of (Preference, Enumeration)

Default: Empty (not setting a preference for any interface)

This property allows the application to select any specific network interfaces or categories of interfaces it wants to Require, Prohibit, Prefer, or Avoid. Note that marking a specific interface as Require strictly limits path selection to that single interface, and often leads to less flexible and resilient connection establishment.

In contrast to other Selection Properties, this property is a tuple of an (Enumerated) interface identifier and a preference, and can either be implemented directly as such, or for making one preference available for each interface and interface type available on the system.

The set of valid interface types is implementation- and system-specific. For example, on a mobile device, there may be Wi-Fi and Cellular interface types available; whereas on a desktop computer, Wi-Fi and Wired Ethernet interface types might be available. An implementation should provide all types that are supported on the local system, to allow applications to be written generically. For example, if a single implementation is used on both mobile devices and desktop devices, it should define the Cellular interface type for both systems, since an application might wish to always prohibit cellular.

The set of interface types is expected to change over time as new access technologies become available. The taxonomy of interface types on a given Transport Services system is implementation-specific.

Interface types should not be treated as a proxy for properties of interfaces such as metered or unmetered network access. If an application needs to prohibit metered interfaces, this should be specified via Provisioning Domain attributes (see Section 6.2.12) or another specific property.

Note that this property is not used to specify an interface scope for a particular endpoint. Section 6.1.2 provides details about how to qualify endpoint candidates on a per-interface basis.

6.2.12. Provisioning Domain Instance or Type

Name: pvd

Type: Collection of (Preference, Enumeration)

Default: Empty (not setting a preference for any PvD)

Similar to interface instances and types (see Section 6.2.11), this property allows the application to control path selection by selecting which specific Provisioning Domain (PvD) or categories of PvDs it wants to Require, Prohibit, Prefer, or Avoid. Provisioning Domains define consistent sets of network properties that may be more specific than network interfaces [RFC7556].

As with interface instances and types, this property is a tuple of an (Enumerated) PvD identifier and a preference, and can either be implemented directly as such, or for making one preference available for each interface and interface type available on the system.

The identification of a specific PvD is implementation- and system-specific, because there is currently no portable standard format for a PvD identifier. For example, this identifier might be a string name or an integer. As with requiring specific interfaces, requiring a specific PvD strictly limits the path selection.

Categories or types of PvDs are also defined to be implementation- and system-specific. These can be useful to identify a service that is provided by a PvD. For example, if an application wants to use a PvD that provides a Voice-Over-IP service on a Cellular network, it can use the relevant PvD type to require a PvD that provides this service, without needing to look up a particular instance. While this does restrict path selection, it is broader than requiring specific PvD instances or interface instances, and should be preferred over these options.

6.2.13. Use Temporary Local Address

Name: useTemporaryLocalAddress

Type: Preference

Default: Avoid for Listeners and Rendezvous Connections. Prefer for other Connections.

This property allows the application to express a preference for the use of temporary local addresses, sometimes called "privacy" addresses [RFC8981]. Temporary addresses are generally used to prevent linking connections over time when a stable address, sometimes called "permanent" address, is not needed. There are some caveats to note when specifying this property. First, if an application Requires the use of temporary addresses, the resulting Connection cannot use IPv4, because temporary addresses do not exist in IPv4. Second, temporary local addresses might involve trading off privacy for performance. For instance, temporary addresses can interfere with resumption mechanisms that some protocols rely on to reduce initial latency.

6.2.14. Multipath Transport

Name: multipath

Type: Enumeration

Default: Disabled for connections created through initiate and rendezvous, Passive for listeners

This property specifies whether and how applications want to take advantage of transferring data across multiple paths between the same end hosts. Using multiple paths allows connections to migrate between interfaces or aggregate bandwidth as availability and performance properties change. Possible values are:

Disabled: The connection will not use multiple paths once established, even if the chosen transport supports using multiple paths.

Active: The connection will negotiate the use of multiple paths if the chosen transport supports this.

Passive: The connection will support the use of multiple paths if the Remote Endpoint requests it.

The policy for using multiple paths is specified using the separate multipath-policy property, see Section 8.1.7 below. To enable the peer endpoint to initiate additional paths towards a local address other than the one initially used, it is necessary to set the Alternative Addresses property (see Section 6.2.15 below).

Setting this property to "Active", can have privacy implications: It enables the transport to establish connectivity using alternate paths that might result in users being linkable across the multiple paths, even if the Advertisement of Alternative Addresses property (see Section 6.2.15 below) is set to false.

Note that Multipath Transport has no corresponding Selection Property of type Preference. Enumeration values other than "Disabled" are interpreted as a preference for choosing protocols that can make use of multiple paths. The "Disabled" value implies a requirement not to use multiple paths in parallel but does not prevent choosing a protocol that is capable of using multiple paths, e.g., it does not prevent choosing TCP, but prevents sending the MP_CAPABLE option in the TCP handshake.

6.2.15. Advertisement of Alternative Addresses

Name: `advertises-altaddr`

Type: Boolean

Default: False

This property specifies whether alternative addresses, e.g., of other interfaces, should be advertised to the peer endpoint by the protocol stack. Advertising these addresses enables the peer-endpoint to establish additional connectivity, e.g., for connection migration or using multiple paths.

Note that this can have privacy implications because it might result in users being linkable across the multiple paths. Also, note that setting this to false does not prevent the local Transport Services system from establishing connectivity using alternate paths (see Section 6.2.14 above); it only prevents proactive advertisement of addresses.

6.2.16. Direction of communication

Name: `direction`

Type: Enumeration

Default: Bidirectional

This property specifies whether an application wants to use the connection for sending and/or receiving data. Possible values are:

Bidirectional: The connection must support sending and receiving

data

Unidirectional send: The connection must support sending data, and the application cannot use the connection to receive any data

Unidirectional receive: The connection must support receiving data, and the application cannot use the connection to send any data

Since unidirectional communication can be supported by transports offering bidirectional communication, specifying unidirectional communication may cause a transport stack that supports bidirectional communication to be selected.

6.2.17. Notification of ICMP soft error message arrival

Name: `softErrorNotify`

Type: Preference

Default: Ignore

This property specifies whether an application considers it useful to be informed when an ICMP error message arrives that does not force termination of a connection. When set to true, received ICMP errors are available as `SoftErrors`, see Section 8.3.1. Note that even if a protocol supporting this property is selected, not all ICMP errors will necessarily be delivered, so applications cannot rely upon receiving them [RFC8085].

6.2.18. Initiating side is not the first to write

Name: `activeReadBeforeSend`

Type: Preference

Default: Ignore

The most common client-server communication pattern involves the client actively opening a connection, then sending data to the server. The server listens (passive open), reads, and then answers. This property specifies whether an application wants to diverge from this pattern - either by actively opening with `Initiate()`, immediately followed by reading, or passively opening with `Listen()`, immediately followed by writing. This property is ignored when establishing connections using `Rendezvous()`. Requiring this property limits the choice of mappings to underlying protocols, which can reduce efficiency. For example, it prevents the Transport Services system from mapping Connections to SCTP streams, where the first transmitted data takes the role of an active open signal [I-D.ietf-taps-impl].

6.3. Specifying Security Parameters and Callbacks

Most security parameters, e.g., TLS ciphersuites, local identity and private key, etc., may be configured statically. Others are dynamically configured during connection establishment. Security parameters and callbacks are partitioned based on their place in the lifetime of connection establishment. Similar to Transport Properties, both parameters and callbacks are inherited during cloning (see Section 7.4).

6.3.1. Specifying Security Parameters on a Pre-Connection

Common security parameters such as TLS ciphersuites are known to implementations. Clients should use common safe defaults for these values whenever possible. However, as discussed in [RFC8922], many transport security protocols require specific security parameters and constraints from the client at the time of configuration and actively during a handshake. These configuration parameters need to be specified in the pre-connection phase and are created as follows:

```
SecurityParameters := NewSecurityParameters()
```

Security configuration parameters and sample usage follow:

- * Local identity and private keys: Used to perform private key operations and prove one's identity to the Remote Endpoint. (Note, if private keys are not available, e.g., since they are stored in hardware security modules (HSMs), handshake callbacks must be used. See below for details.)

```
SecurityParameters.Set(identity, myIdentity)  
SecurityParameters.Set(key-pair, myPrivateKey, myPublicKey)
```

- * Supported algorithms: Used to restrict what parameters are used by underlying transport security protocols. When not specified, these algorithms should use known and safe defaults for the system. Parameters include: ciphersuites, supported groups, and signature algorithms. These parameters take a collection of supported algorithms as parameter.

```
SecurityParameters.Set(supported-group, "secp256r1")
SecurityParameters.Set(ciphersuite, "TLS_AES_128_GCM_SHA256")
SecurityParameters.Set(signature-algorithm, "ecdsa_secp256r1_sha256")
```

- * Pre-Shared Key import: Used to install pre-shared keying material established out-of-band. Each pre-shared keying material is associated with some identity that typically identifies its use or has some protocol-specific meaning to the Remote Endpoint.

```
SecurityParameters.Set(pre-shared-key, key, identity)
```

- * Session cache management: Used to tune session cache capacity, lifetime, and other policies.

```
SecurityParameters.Set(max-cached-sessions, 16)
SecurityParameters.Set(cached-session-lifetime-seconds, 3600)
```

Connections that use Transport Services SHOULD use security in general. However, for compatibility with endpoints that do not support transport security protocols (such as a TCP endpoint that does not support TLS), applications can initialize their security parameters to indicate that security can be disabled, or can be opportunistic. If security is disabled, the Transport Services system will not attempt to add transport security automatically. If security is opportunistic, it will allow Connections without transport security, but will still attempt to use security if available.

```
SecurityParameters := NewDisabledSecurityParameters()
```

```
SecurityParameters := NewOpportunisticSecurityParameters()
```

Representation of Security Parameters in implementations should parallel that chosen for Transport Property names as suggested in Section 5.

6.3.2. Connection Establishment Callbacks

Security decisions, especially pertaining to trust, are not static. Once configured, parameters may also be supplied during connection establishment. These are best handled as client-provided callbacks. Callbacks block the progress of the connection establishment, which distinguishes them from other Events in the transport system. How callbacks and events are implemented is specific to each implementation. Security handshake callbacks that may be invoked during connection establishment include:

- * Trust verification callback: Invoked when a Remote Endpoint's trust must be verified before the handshake protocol can continue. For example, the application could verify an X.509 certificate as described in [RFC5280].

```
TrustCallback := NewCallback({  
    // Handle trust, return the result  
})  
SecurityParameters.SetTrustVerificationCallback(trustCallback)
```

- * Identity challenge callback: Invoked when a private key operation is required, e.g., when local authentication is requested by a Remote Endpoint.

```
ChallengeCallback := NewCallback({  
    // Handle challenge  
})  
SecurityParameters.SetIdentityChallengeCallback(challengeCallback)
```

7. Establishing Connections

Before a Connection can be used for data transfer, it needs to be established. Establishment ends the pre-establishment phase; all transport properties and cryptographic parameter specification must be complete before establishment, as these will be used to select candidate Paths and Protocol Stacks for the Connection. Establishment may be active, using the Initiate() Action; passive, using the Listen() Action; or simultaneous for peer-to-peer, using the Rendezvous() Action. These Actions are described in the subsections below.

7.1. Active Open: Initiate

Active open is the Action of establishing a Connection to a Remote Endpoint presumed to be listening for incoming Connection requests. Active open is used by clients in client-server interactions. Active open is supported by the Transport Services API through the Initiate Action:

```
Connection := Preconnection.Initiate(timeout?)
```

The timeout parameter specifies how long to wait before aborting Active open. Before calling Initiate, the caller must have populated a Preconnection Object with a Remote Endpoint specifier, optionally a Local Endpoint specifier (if not specified, the system will attempt to determine a suitable Local Endpoint), as well as all properties necessary for candidate selection.

The Initiate() Action returns a Connection object. Once Initiate() has been called, any changes to the Preconnection MUST NOT have any effect on the Connection. However, the Preconnection can be reused, e.g., to Initiate another Connection.

Once Initiate is called, the candidate Protocol Stack(s) may cause one or more candidate transport-layer connections to be created to the specified Remote Endpoint. The caller may immediately begin sending Messages on the Connection (see Section 9.2) after calling Initiate(); note that any data marked Safely Replayable that is sent while the Connection is being established may be sent multiple times or on multiple candidates.

The following Events may be sent by the Connection after Initiate() is called:

```
Connection -> Ready<>
```

The Ready Event occurs after Initiate has established a transport-layer connection on at least one usable candidate Protocol Stack over at least one candidate Path. No Receive Events (see Section 9.3) will occur before the Ready Event for Connections established using Initiate.

```
Connection -> EstablishmentError<reason?>
```

An EstablishmentError occurs either when the set of transport properties and security parameters cannot be fulfilled on a Connection for initiation (e.g., the set of available Paths and/or Protocol Stacks meeting the constraints is empty) or reconciled with the Local and/or Remote Endpoints; when the remote specifier cannot

be resolved; or when no transport-layer connection can be established to the Remote Endpoint (e.g., because the Remote Endpoint is not accepting connections, the application is prohibited from opening a Connection by the operating system, or the establishment attempt has timed out for any other reason).

Connection establishment and transmission of the first message can be combined in a single action Section 9.2.5.

7.2. Passive Open: Listen

Passive open is the Action of waiting for Connections from Remote Endpoints, commonly used by servers in client-server interactions. Passive open is supported by the Transport Services API through the Listen Action and returns a Listener object:

```
Listener := Preconnection.Listen()
```

Before calling Listen, the caller must have initialized the Preconnection during the pre-establishment phase with a Local Endpoint specifier, as well as all properties necessary for Protocol Stack selection. A Remote Endpoint may optionally be specified, to constrain what Connections are accepted.

The Listen() Action returns a Listener object. Once Listen() has been called, any changes to the Preconnection MUST NOT have any effect on the Listener. The Preconnection can be disposed of or reused, e.g., to create another Listener.

```
Listener.Stop()
```

Listening continues until the global context shuts down, or until the Stop action is performed on the Listener object.

```
Listener -> ConnectionReceived<Connection>
```

The ConnectionReceived Event occurs when a Remote Endpoint has established a transport-layer connection to this Listener (for Connection-oriented transport protocols), or when the first Message has been received from the Remote Endpoint (for Connectionless protocols), causing a new Connection to be created. The resulting Connection is contained within the ConnectionReceived Event, and is ready to use as soon as it is passed to the application via the event.

```
Listener.SetNewConnectionLimit(value)
```

If the caller wants to rate-limit the number of inbound Connections that will be delivered, it can set a cap using `SetNewConnectionLimit()`. This mechanism allows a server to protect itself from being drained of resources. Each time a new Connection is delivered by the `ConnectionReceived` Event, the value is automatically decremented. Once the value reaches zero, no further Connections will be delivered until the caller sets the limit to a higher value. By default, this value is Infinite. The caller is also able to reset the value to Infinite at any point.

Listener -> EstablishmentError<reason?>

An EstablishmentError occurs either when the Properties and Security Parameters of the Preconnection cannot be fulfilled for listening or cannot be reconciled with the Local Endpoint (and/or Remote Endpoint, if specified), when the Local Endpoint (or Remote Endpoint, if specified) cannot be resolved, or when the application is prohibited from listening by policy.

Listener -> Stopped<>

A Stopped Event occurs after the Listener has stopped listening.

7.3. Peer-to-Peer Establishment: Rendezvous

Simultaneous peer-to-peer Connection establishment is supported by the `Rendezvous()` Action:

`Preconnection.Rendezvous()`

A Preconnection Object used in a `Rendezvous()` MUST have both the Local Endpoint candidates and the Remote Endpoint candidates specified, along with the transport properties and security parameters needed for Protocol Stack selection, before the `Rendezvous()` Action is initiated.

The `Rendezvous()` Action listens on the Local Endpoint candidates for an incoming Connection from the Remote Endpoint candidates, while also simultaneously trying to establish a Connection from the Local Endpoint candidates to the Remote Endpoint candidates.

If there are multiple Local Endpoints or Remote Endpoints configured, then initiating a `Rendezvous()` action will systematically probe the reachability of those endpoint candidates following an approach such as that used in Interactive Connectivity Establishment (ICE) [RFC8445].

If the endpoints are suspected to be behind a NAT, `Rendezvous()` can be initiated using Local Endpoints that support a method of discovering NAT bindings such as Session Traversal Utilities for NAT (STUN) [RFC8489] or Traversal Using Relays around NAT (TURN) [RFC8656]. In this case, the Local Endpoint will resolve to a mixture of local and server reflexive addresses. The `Resolve()` action on the Preconnection can be used to discover these bindings:

```
[]LocalEndpoint, []RemoteEndpoint := Preconnection.Resolve()
```

The `Resolve()` call returns lists of Local Endpoints and Remote Endpoints, that represent the concrete addresses, local and server reflexive, on which a `Rendezvous()` for the Preconnection will listen for incoming Connections, and to which it will attempt to establish connections.

Note that the set of LocalEndpoints returned by `Resolve()` might or might not contain information about all possible local interfaces; it is valid only for a `Rendezvous` happening at the same time as the resolution. Care should be taken in using these values in any other context.

An application that uses `Rendezvous()` to establish a peer-to-peer connection in the presence of NATs will configure the Preconnection object with at least one a Local Endpoint that supports NAT binding discovery. It will then `Resolve()` the Preconnection, and pass the resulting list of Local Endpoint candidates to the peer via a signalling protocol, for example as part of an ICE [RFC5245] exchange within SIP [RFC3261] or WebRTC [RFC7478]. The peer will then, via the same signalling channel, return the Remote Endpoint candidates. The set of Remote Endpoint candidates are then configured onto the Preconnection:

```
Preconnection.AddRemote([]RemoteEndpoint)
```

The `Rendezvous()` Action can be initiated once both the Local Endpoint candidates and the Remote Endpoint candidates retrieved from the peer via the signalling channel have been added to the Preconnection.

If successful, the `Rendezvous()` Action returns a Connection object via a `RendezvousDone<>` Event:

```
Preconnection -> RendezvousDone<Connection>
```

The `RendezvousDone<>` Event occurs when a Connection is established with the Remote Endpoint. For Connection-oriented transports, this occurs when the transport-layer connection is established; for Connectionless transports, it occurs when the first Message is

received from the Remote Endpoint. The resulting Connection is contained within the RendezvousDone<> Event, and is ready to use as soon as it is passed to the application via the Event. Changes made to a Preconnection after Rendezvous() has been called do not have any effect on existing Connections.

An EstablishmentError occurs either when the Properties and Security Parameters of the Preconnection cannot be fulfilled for rendezvous or cannot be reconciled with the Local and/or Remote Endpoints, when the Local Endpoint or Remote Endpoint cannot be resolved, when no transport-layer connection can be established to the Remote Endpoint, or when the application is prohibited from rendezvous by policy:

Preconnection -> EstablishmentError<reason?>

7.4. Connection Groups

Connection Groups can be created using the Clone Action:

Connection := Connection.Clone(framer?)

Calling Clone on a Connection yields a Connection Group containing two Connections: the parent Connection on which Clone was called, and a resulting cloned Connection. The new Connection is actively opened, and it will send a Ready Event or an EstablishmentError Event. Calling Clone on any of these Connections adds another Connection to the Connection Group. Connections in a Connection Group share all Connection Properties except Connection Priority (see Section 8.1.2), and these Connection Properties are entangled: Changing one of the Connection Properties on one Connection in the Connection Group automatically changes the Connection Property for all others. For example, changing Timeout for aborting Connection (see Section 8.1.3) on one Connection in a Connection Group will automatically make the same change to this Connection Property for all other Connections in the Connection Group. Like all other Properties, Connection Priority is copied to the new Connection when calling Clone(), but in this case, a later change to the Connection Priority on one Connection does not change it on the other Connections in the same Connection Group.

Message Properties set on a Connection also apply only to that Connection.

A new Connection created by Clone can have a Message Framers assigned via the optional framer parameter of the Clone Action. If this parameter is not supplied, the stack of Message Framers associated with a Connection is copied to the cloned Connection when calling Clone. Then, a cloned Connection has the same stack of Message Framers as the Connection from which they are Cloned, but these Framers may internally maintain per-Connection state.

It is also possible to check which Connections belong to the same Connection Group. Calling GroupedConnections() on a specific Connection returns a set of all Connections in the same group.

```
[]Connection := Connection.GroupedConnections()
```

Connections will belong to the same group if the application previously called Clone. Passive Connections can also be added to the same group - e.g., when a Listener receives a new Connection that is just a new stream of an already active multi-streaming protocol instance.

If the underlying protocol supports multi-streaming, it is natural to use this functionality to implement Clone. In that case, Connections in a Connection Group are multiplexed together, giving them similar treatment not only inside endpoints, but also across the end-to-end Internet path.

Note that calling Clone() can result in on-the-wire signaling, e.g., to open a new transport connection, depending on the underlying Protocol Stack. When Clone() leads to the opening of multiple such connections, the Transport Services system will ensure consistency of Connection Properties by uniformly applying them to all underlying connections in a group. Even in such a case, there are possibilities for a Transport Services system to implement prioritization within a Connection Group [TCP-COUPLING] [RFC8699].

Attempts to clone a Connection can result in a CloneError:

```
Connection -> CloneError<reason?>
```

The Connection Priority Connection Property operates on Connections in a Connection Group using the same approach as in Section 9.1.3.2: when allocating available network capacity among Connections in a Connection Group, sends on Connections with higher Priority values will be prioritized over sends on Connections that have lower Priority values. Capacity will be shared among these Connections according to the Connection Group Transmission Scheduler property (Section 8.1.5). See Section 9.2.6 for more.

7.5. Adding and Removing Endpoints on a Connection

Transport protocols that are explicitly multipath aware are expected to automatically manage the set of Remote Endpoints that they are communicating with, and the paths to those endpoints. A PathChange<> event, described in Section 8.3.2, will be generated when the path changes.

In some cases, however, it is necessary to explicitly indicate to a Connection that a new remote endpoint has become available for use, or to indicate that some remote endpoint is no longer available. This is most common in the case of peer to peer connections using Trickle ICE [RFC8838].

The AddRemote() action can be used to add one or more new remote endpoints to a Connection:

```
Connection.AddRemote([]RemoteEndpoint)
```

Endpoints that are already known to the Connection are ignored. A call to AddRemote() makes the new remote endpoints available to the connection, but whether the Connection makes use of those endpoints will depend on the underlying transport protocol.

Similarly, the RemoveRemote() action can be used to tell a connection to stop using one or more remote endpoints:

```
Connection.RemoveRemote([]RemoteEndpoint)
```

Removing all known remote endpoints can have the effect of aborting the connection. The effect of removing the active remote endpoint(s) depends on the underlying transport: multipath aware transports might be able to switch to a new path if other reachable remote endpoints exist, or the connection might abort.

Similarly, the AddLocal() and RemoveLocal() actions can be used to add and remove local endpoints to/from a Connection.

8. Managing Connections

During pre-establishment and after establishment, connections can be configured and queried using Connection Properties, and asynchronous information may be available about the state of the connection via Soft Errors.

Connection Properties represent the configuration and state of the selected Protocol Stack(s) backing a Connection. These Connection Properties may be Generic, applying regardless of transport protocol,

or Specific, applicable to a single implementation of a single transport protocol stack. Generic Connection Properties are defined in Section 8.1 below.

Protocol Specific Properties are defined in a transport- and implementation-specific way to permit more specialized protocol features to be used. Too much reliance by an application on Protocol Specific Properties can significantly reduce the flexibility of a transport services implementation to make appropriate selection and configuration choices. Therefore, it is RECOMMENDED that Protocol Properties are used for properties common across different protocols and that Protocol Specific Properties are only used where specific protocols or properties are necessary.

The application can set and query Connection Properties on a per-Connection basis. Connection Properties that are not read-only can be set during pre-establishment (see Section 6.2), as well as on connections directly using the SetProperty action:

```
Connection.SetProperty(property, value)
```

Note that changing one of the Connection Properties on one Connection in a Connection Group will also change it for all other Connections of that group; see further Section 7.4.

At any point, the application can query Connection Properties.

```
ConnectionProperties := Connection.GetProperties()  
value := ConnectionProperties.Get(property)  
if ConnectionProperties.Has(boolean_or_preference_property) then ...
```

Depending on the status of the connection, the queried Connection Properties will include different information:

- * The connection state, which can be one of the following:
Establishing, Established, Closing, or Closed.
- * Whether the connection can be used to send data. A connection can not be used for sending if the connection was created with the Selection Property Direction of Communication set to unidirectional receive or if a Message marked as Final was sent over this connection. See also Section 9.1.3.5.

- * Whether the connection can be used to receive data. A connection cannot be used for reading if the connection was created with the Selection Property Direction of Communication set to unidirectional send or if a Message marked as Final was received. See Section 9.3.3.3. The latter is only supported by certain transport protocols, e.g., by TCP as half-closed connection.
- * For Connections that are Established, Closing, or Closed: Connection Properties (Section 8.1) of the actual protocols that were selected and instantiated, and Selection Properties that the application specified on the Preconnection. Selection Properties of type Preference will be exposed as boolean values indicating whether or not the property applies to the selected transport. Note that the instantiated protocol stack might not match all Protocol Selection Properties that the application specified on the Preconnection.
- * For Connections that are Established: information concerning the path(s) used by the Protocol Stack. This can be derived from local PVD information, measurements by the Protocol Stack, or other sources. For example, a TAPS system that is configured to receive and process PVD information [RFC7556] could also provide network configuration information for the chosen path(s).

8.1. Generic Connection Properties

Generic Connection Properties are defined independent of the chosen protocol stack and therefore available on all Connections.

Many Connection Properties have a corresponding Selection Property that enables applications to express their preference for protocols providing a supporting transport feature.

8.1.1. Required Minimum Corruption Protection Coverage for Receiving

Name: `recvChecksumLen`

Type: Integer or Full Coverage

Default: Full Coverage

If this property is an Integer, it specifies the minimum number of bytes in a received message that need to be covered by a checksum. A receiving endpoint will not forward messages that have less coverage to the application. The application is responsible for handling any corruption within the non-protected part of the message [RFC8085]. A special value of 0 means that a received packet may also have a zero checksum field.

8.1.2. Connection Priority

Name: connPriority

Type: Integer (non-negative)

Default: 100

This Property is a non-negative integer representing the priority of this Connection relative to other Connections in the same Connection Group. A higher value reflects a higher priority. It has no effect on Connections not part of a Connection Group. As noted in Section 7.4, this property is not entangled when Connections are cloned, i.e., changing the Priority on one Connection in a Connection Group does not change it on the other Connections in the same Connection Group. No guarantees of a specific behavior regarding Connection Priority are given; a Transport Services system may ignore this property. See Section 9.2.6 for more details.

8.1.3. Timeout for Aborting Connection

Name: connTimeout

Type: Numeric or Disabled

Default: Disabled

If this property is Numeric, it specifies how long to wait before deciding that an active Connection has failed when trying to reliably deliver data to the Remote Endpoint. Adjusting this Property will only take effect when the underlying stack supports reliability. If this property has the enumerated value Disabled, it means that no timeout is scheduled.

8.1.4. Timeout for keep alive packets

Name: keepAliveTimeout

Type: Numeric or Disabled

Default: Implementation-defined

A Transport Services API can request a protocol that supports sending keep alive packets Section 6.2.10. If this property is an Integer, it specifies the maximum length of time an idle connection (one for which no transport packets have been sent) should wait before the Local Endpoint sends a keep-alive packet to the Remote Endpoint. Adjusting this Property will only take effect when the underlying

stack supports sending keep-alive packets. Guidance on setting this value for datagram transports is provided in [RFC8085]. A value greater than the connection timeout (Section 8.1.3) or the enumerated value Disabled will disable the sending of keep-alive packets.

8.1.5. Connection Group Transmission Scheduler

Name: connScheduler

Type: Enumeration

Default: Weighted Fair Queueing (see Section 3.6 in [RFC8260])

This property specifies which scheduler should be used among Connections within a Connection Group, see Section 7.4. The set of schedulers can be taken from [RFC8260].

8.1.6. Capacity Profile

Name: connCapacityProfile

Type: Enumeration

Default: Default Profile (Best Effort)

This property specifies the desired network treatment for traffic sent by the application and the tradeoffs the application is prepared to make in path and protocol selection to receive that desired treatment. When the capacity profile is set to a value other than Default, z Transport Services system SHOULD select paths and configure protocols to optimize the tradeoff between delay, delay variation, and efficient use of the available capacity based on the capacity profile specified. How this is realized is implementation-specific. The Capacity Profile MAY also be used to set markings on the wire for Protocol Stacks supporting this. Recommendations for use with DSCP are provided below for each profile; note that when a Connection is multiplexed, the guidelines in Section 6 of [RFC7657] apply.

The following values are valid for the Capacity Profile:

Default: The application provides no information about its expected capacity profile. Transport Services implementations that map the requested capacity profile onto per-connection DSCP signaling SHOULD assign the DSCP Default Forwarding [RFC2474] Per Hop Behaviour (PHB).

Scavenger: The application is not interactive. It expects to send

and/or receive data without any urgency. This can, for example, be used to select protocol stacks with scavenger transmission control and/or to assign the traffic to a lower-effort service. Transport Services implementations that map the requested capacity profile onto per-connection DSCP signaling SHOULD assign the DSCP Less than Best Effort [RFC8622] PHB.

Low Latency/Interactive: The application is interactive, and prefers loss to latency. Response time should be optimized at the expense of delay variation and efficient use of the available capacity when sending on this connection. This can be used by the system to disable the coalescing of multiple small Messages into larger packets (Nagle's algorithm); to prefer immediate acknowledgment from the peer endpoint when supported by the underlying transport; and so on. Transport Services implementations that map the requested capacity profile onto per-connection DSCP signaling without multiplexing SHOULD assign a DSCP Assured Forwarding (AF41,AF42,AF43,AF44) [RFC2597] PHB. Inelastic traffic that is expected to conform to the configured network service rate could be mapped to the DSCP Expedited Forwarding [RFC3246] or [RFC5865] PHBs.

Low Latency/Non-Interactive: The application prefers loss to latency, but is not interactive. Response time should be optimized at the expense of delay variation and efficient use of the available capacity when sending on this connection. Transport system implementations that map the requested capacity profile onto per-connection DSCP signaling without multiplexing SHOULD assign a DSCP Assured Forwarding (AF21,AF22,AF23,AF24) [RFC2597] PHB.

Constant-Rate Streaming: The application expects to send/receive data at a constant rate after Connection establishment. Delay and delay variation should be minimized at the expense of efficient use of the available capacity. This implies that the Connection might fail if the Path is unable to maintain the desired rate. A transport can interpret this capacity profile as preferring a circuit breaker [RFC8084] to a rate-adaptive congestion controller. Transport system implementations that map the requested capacity profile onto per-connection DSCP signaling without multiplexing SHOULD assign a DSCP Assured Forwarding (AF31,AF32,AF33,AF34) [RFC2597] PHB.

Capacity-Seeking: The application expects to send/receive data at

the maximum rate allowed by its congestion controller over a relatively long period of time. Transport Services implementations that map the requested capacity profile onto per-connection DSCP signaling without multiplexing SHOULD assign a DSCP Assured Forwarding (AF11,AF12,AF13,AF14) [RFC2597] PHB per Section 4.8 of [RFC4594].

The Capacity Profile for a selected protocol stack may be modified on a per-Message basis using the Transmission Profile Message Property; see Section 9.1.3.8.

8.1.7. Policy for using Multipath Transports

Name: multipath-policy

Type: Enumeration

Default: Handover

This property specifies the local policy for transferring data across multiple paths between the same end hosts if Multipath Transport is not set to Disabled (see Section 6.2.14). Possible values are:

Handover: The connection ought only to attempt to migrate between different paths when the original path is lost or becomes unusable. The thresholds used to declare a path unusable are implementation specific.

Interactive: The connection ought only to attempt to minimize the latency for interactive traffic patterns by transmitting data across multiple paths when this is beneficial. The goal of minimizing the latency will be balanced against the cost of each of these paths. Depending on the cost of the lower-latency path, the scheduling might choose to use a higher-latency path. Traffic can be scheduled such that data may be transmitted on multiple paths in parallel to achieve a lower latency. The specific scheduling algorithm is implementation-specific.

Aggregate: The connection ought to attempt to use multiple paths in parallel to maximize available capacity and possibly overcome the capacity limitations of the individual paths. The actual strategy is implementation specific.

Note that this is a local choice - the Remote Endpoint can choose a different policy.

8.1.8. Bounds on Send or Receive Rate

Name: minSendRate / minRecvRate / maxSendRate / maxRecvRate

Type: Numeric or Unlimited / Numeric or Unlimited / Numeric or Unlimited / Numeric or Unlimited

Default: Unlimited / Unlimited / Unlimited / Unlimited

Integer values of this property specify an upper-bound rate that a transfer is not expected to exceed (even if flow control and congestion control allow higher rates), and/or a lower-bound rate below which the application does not deem it will be useful. These are specified in bits per second. The enumerated value Unlimited indicates that no bound is specified.

8.1.9. Group Connection Limit

Name: groupConnLimit

Type: Numeric or Unlimited

Default: Unlimited

If this property is an Integer, it controls the number of Connections that can be accepted from a peer as new members of the Connection's group. Similar to SetNewConnectionLimit(), this limits the number of ConnectionReceived Events that will occur, but constrained to the group of the Connection associated with this property. For a multi-streaming transport, this limits the number of allowed streams.

8.1.10. Isolate Session

Name: isolateSession

Type: Boolean

Default: false

When set to true, this property will initiate new Connections using as little cached information (such as session tickets or cookies) as possible from previous connections that are not in the same Connection Group. Any state generated by this Connection will only be shared with Connections in the same Connection Group. Cloned Connections will use saved state from within the Connection Group. This is used for separating Connection Contexts as specified in [I-D.ietf-taps-arch].

Note that this does not guarantee no leakage of information, as implementations may not be able to fully isolate all caches (e.g. RTT estimates). Note that this property may degrade connection performance.

8.1.11. Read-only Connection Properties

The following generic Connection Properties are read-only, i.e. they cannot be changed by an application.

8.1.11.1. Maximum Message Size Concurrent with Connection Establishment

Name: zeroRttMsgMaxLen

Type: Integer

This property represents the maximum Message size that can be sent before or during Connection establishment, see also Section 9.1.3.4. It is specified as the number of bytes.

8.1.11.2. Maximum Message Size Before Fragmentation or Segmentation

Name: singularTransmissionMsgMaxLen

Type: Integer

This property, if applicable, represents the maximum Message size that can be sent without incurring network-layer fragmentation at the sender. It is specified as the number of bytes. It exposes a value to the application based on the Maximum Packet Size (MPS) as described in Datagram PLPMTUD [RFC8899]. This can allow a sending stack to avoid unwanted fragmentation at the network-layer or segmentation by the transport layer.

8.1.11.3. Maximum Message Size on Send

Name: sendMsgMaxLen

Type: Integer

This property represents the maximum Message size that an application can send. It is specified as the number of bytes.

8.1.11.4. Maximum Message Size on Receive

Name: recvMsgMaxLen

Type: Integer

This numeric property represents the maximum Message size that an application can receive. It specified as the number of bytes.

8.2. TCP-specific Properties: User Timeout Option (UTO)

These properties specify configurations for the User Timeout Option (UTO), in the case that TCP becomes the chosen transport protocol. Implementation is optional and useful only if TCP is implemented in the Transport Services system.

These TCP-specific properties are included here because the feature Suggest timeout to the peer is part of the minimal set of transport services [RFC8923], where this feature was categorized as "functional". This means that when an Transport Services implementation offers this feature, the Transport Services API has to expose an interface to the application. Otherwise, the implementation might violate assumptions by the application, which could cause the application to fail.

All of the below properties are optional (e.g., it is possible to specify User Timeout Enabled as true, but not specify an Advertised User Timeout value; in this case, the TCP default will be used). These properties reflect the API extension specified in Section 3 of [RFC5482].

8.2.1. Advertised User Timeout

Name: tcp.userTimeoutValue

Type: Integer

Default: the TCP default

This time value is advertised via the TCP User Timeout Option (UTO) [RFC5482] at the Remote Endpoint to adapt its own Timeout for aborting Connection (see Section 8.1.3) value.

8.2.2. User Timeout Enabled

Name: tcp.userTimeoutEnabled

Type: Boolean

Default: false

This property controls whether the UTO option is enabled for a connection. This applies to both sending and receiving.

8.2.3. Timeout Changeable

Name: tcp.userTimeoutChangeable

Type: Boolean

Default: true

This property controls whether the Timeout for aborting Connection (see Section 8.1.3) may be changed based on a UTO option received from the remote peer. This boolean becomes false when Timeout for aborting Connection (see Section 8.1.3) is used.

8.3. Connection Lifecycle Events

During the lifetime of a connection there are events that can occur when configured.

8.3.1. Soft Errors

Asynchronous introspection is also possible, via the SoftError Event. This event informs the application about the receipt and contents of an ICMP error message related to the Connection. This will only happen if the underlying protocol stack supports access to soft errors; however, even if the underlying stack supports it, there is no guarantee that a soft error will be signaled.

Connection -> SoftError<>

8.3.2. Path change

This event notifies the application when at least one of the paths underlying a Connection has changed. Changes occur on a single path when the PMTU changes as well as when multiple paths are used and paths are added or removed, the set of local endpoints changes, or a handover has been performed.

Connection -> PathChange<>

9. Data Transfer

Data is sent and received as Messages, which allows the application to communicate the boundaries of the data being transferred.

9.1. Messages and Framers

Each Message has an optional Message Context, which allows to add Message Properties, identify Send Events related to a specific Message or to inspect meta-data related to the Message sent. Framers can be used to extend or modify the message data with additional information that can be processed at the receiver to detect message boundaries.

9.1.1. Message Contexts

Using the MessageContext object, the application can set and retrieve meta-data of the message, including Message Properties (see Section 9.1.3) and framing meta-data (see Section 9.1.2.2). Therefore, a MessageContext object can be passed to the Send action and is returned by each Send and Receive related event.

Message Properties can be set and queried using the Message Context:

```
MessageContext.add(property, value)
PropertyValue := MessageContext.get(property)
```

These Message Properties may be generic properties or Protocol Specific Properties.

For MessageContexts returned by send Events (see Section 9.2.2) and receive Events (see Section 9.3.2), the application can query information about the Local and Remote Endpoint:

```
RemoteEndpoint := MessageContext.GetRemoteEndpoint()
LocalEndpoint := MessageContext.GetLocalEndpoint()
```

9.1.2. Message Framers

Although most applications communicate over a network using well-formed Messages, the boundaries and metadata of the Messages are often not directly communicated by the transport protocol itself. For example, HTTP applications send and receive HTTP messages over a byte-stream transport, requiring that the boundaries of HTTP messages be parsed from the stream of bytes.

Message Framers allow extending a Connection's Protocol Stack to define how to encapsulate or encode outbound Messages, and how to decapsulate or decode inbound data into Messages. Message Framers allow message boundaries to be preserved when using a Connection object, even when using byte-stream transports. This is designed based on the fact that many of the current application protocols evolved over TCP, which does not provide message boundary

preservation, and since many of these protocols require message boundaries to function, each application layer protocol has defined its own framing.

To use a Message Framer, the application adds it to its Preconnection object. Then, the Message Framer can intercept all calls to `Send()` or `Receive()` on a Connection to add Message semantics, in addition to interacting with the setup and teardown of the Connection. A Framer can start sending data before the application sends data if the framing protocol requires a prefix or handshake (see [RFC8229] for an example of such a framing protocol).

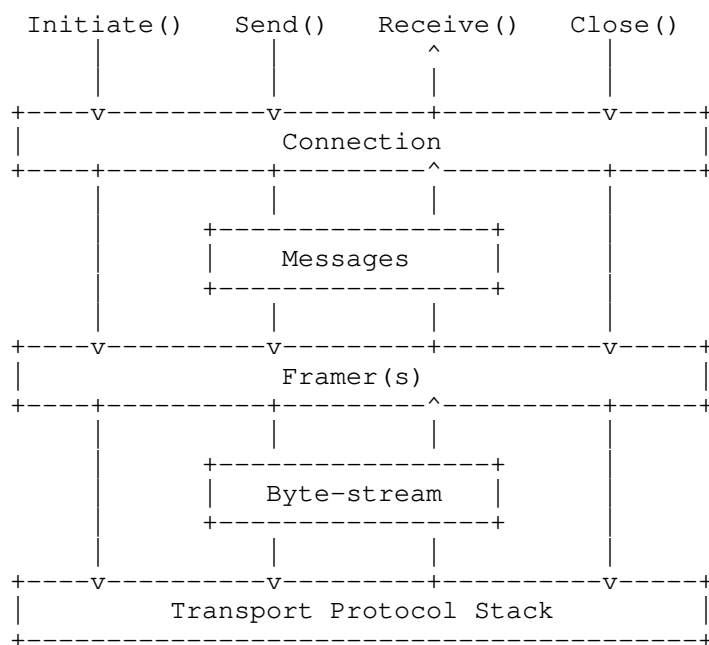


Figure 1: Protocol Stack showing a Message Framer

Note that while Message Framers add the most value when placed above a protocol that otherwise does not preserve message boundaries, they can also be used with datagram- or message-based protocols. In these cases, they add an additional transformation to further encode or encapsulate, and can potentially support packing multiple application-layer Messages into individual transport datagrams.

The API to implement a Message Framer can vary depending on the implementation; guidance on implementing Message Framers can be found in [I-D.ietf-taps-impl].

9.1.2.1. Adding Message Framers to Pre-Connections

The Message Framer object can be added to one or more Preconnections to run on top of transport protocols. Multiple Framers may be added to a Preconnection; in this case, the Framers operate as a framing stack, i.e. the last one added runs first when framing outbound messages, and last when parsing inbound data.

The following example adds a basic HTTP Message Framer to a Preconnection:

```
framer := NewHTTPMessageFramer()  
Preconnection.AddFramer(framer)
```

Since Message Framers pass from Preconnection to Listener or Connection, addition of Framers must happen before any operation that may result in the creation of a Connection.

9.1.2.2. Framing Meta-Data

When sending Messages, applications can add Framer-specific properties to a MessageContext (Section 9.1.1). In order to set these properties, the add and get actions on the MessageContext. To avoid naming conflicts, the property names SHOULD be prefixed with a namespace referencing the framer implementation or the protocol it implements as described in Section 4.1.

This mechanism can be used, for example, to set the type of a Message for a TLV format. The namespace of values is custom for each unique Message Framer.

```
messageContext := NewMessageContext()  
messageContext.add(framer, key, value)  
Connection.Send(messageData, messageContext)
```

When an application receives a MessageContext in a Receive event, it can also look to see if a value was set by a specific Message Framer.

```
messageContext.get(framer, key) -> value
```

For example, if an HTTP Message Framer is used, the values could correspond to HTTP headers:

```
httpFramer := NewHTTPMessageFramer()  
...  
messageContext := NewMessageContext()  
messageContext.add(httpFramer, "accept", "text/html")
```

9.1.3. Message Properties

Applications needing to annotate the Messages they send with extra information (for example, to control how data is scheduled and processed by the transport protocols supporting the Connection) can include this information in the Message Context passed to the Send Action. For other uses of the message context, see Section 9.1.1.

Message Properties are per-Message, not per-Send if partial Messages are sent (Section 9.2.3). All data blocks associated with a single Message share properties specified in the Message Contexts. For example, it would not make sense to have the beginning of a Message expire, but allow the end of a Message to still be sent.

A MessageContext object contains metadata for the Messages to be sent or received.

```
messageData := "hello"
messageContext := NewMessageContext()
messageContext.add(parameter, value)
Connection.Send(messageData, messageContext)
```

The simpler form of Send, which does not take any messageContext, is equivalent to passing a default MessageContext without adding any Message Properties.

If an application wants to override Message Properties for a specific message, it can acquire an empty MessageContext Object and add all desired Message Properties to that Object. It can then reuse the same messageContext Object for sending multiple Messages with the same properties.

Properties can be added to a MessageContext object only before the context is used for sending. Once a MessageContext has been used with a Send call, further modifications to the MessageContext object do not have any effect on this Send call. Message Properties that are not added to a MessageContext object before using the context for sending will either take a specific default value or be configured based on Selection or Connection Properties of the Connection that is associated with the Send call. This initialization behavior is defined per Message Property below.

The Message Properties could be inconsistent with the properties of the Protocol Stacks underlying the Connection on which a given Message is sent. For example, a Protocol Stack must be able to provide ordering if the msgOrdered property of a Message is enabled. Sending a Message with Message Properties inconsistent with the Selection Properties of the Connection yields an error.

If a Message Property contradicts a Connection Property, and if this per-Message behavior can be supported, it overrides the Connection Property for the specific Message. For example, if Reliable Data Transfer (Connection) is set to Require and a protocol with configurable per-Message reliability is used, setting Reliable Data Transfer (Message) to false for a particular Message will allow this Message to be sent without any reliability guarantees. Changing the Reliable Data Transfer property on Messages is only possible for Connections that were established enabling the Selection Property Configure Per-Message Reliability.

The following Message Properties are supported:

9.1.3.1. Lifetime

Name: msgLifetime

Type: Numeric

Default: infinite

The Lifetime specifies how long a particular Message can wait to be sent to the Remote Endpoint before it is irrelevant and no longer needs to be (re-)transmitted. This is a hint to the Transport Services implementation - it is not guaranteed that a Message will not be sent when its Lifetime has expired.

Setting a Message's Lifetime to infinite indicates that the application does not wish to apply a time constraint on the transmission of the Message, but it does not express a need for reliable delivery; reliability is adjustable per Message via the Reliable Data Transfer (Message) property (see Section 9.1.3.7). The type and units of Lifetime are implementation-specific.

9.1.3.2. Priority

Name: msgPriority

Type: Integer (non-negative)

Default: 100

This property specifies the priority of a Message, relative to other Messages sent over the same Connection.

A Message with Priority 0 will yield to a Message with Priority 1, which will yield to a Message with Priority 2, and so on. Priorities may be used as a sender-side scheduling construct only, or be used to specify priorities on the wire for Protocol Stacks supporting prioritization.

Note that this property is not a per-message override of the Connection Priority – see Section 8.1.2. The Priority properties may interact, but can be used independently and be realized by different mechanisms; see Section 9.2.6.

9.1.3.3. Ordered

Name: msgOrdered

Type: Boolean

Default: the queried Boolean value of the Selection Property preserveOrder (Section 6.2.4)

The order in which Messages were submitted for transmission via the Send Action will be preserved on delivery via Receive<> events for all Messages on a Connection that have this Message Property set to true.

If false, the Message is delivered to the receiving application without preserving the ordering. This property is used for protocols that support preservation of data ordering, see Section 6.2.4, but allow out-of-order delivery for certain messages, e.g., by multiplexing independent messages onto different streams.

If it is not configured by the application before sending, this property's default value will be based on the Selection Property preserveOrder of the Connection associated with the Send Action.

9.1.3.4. Safely Replayable

Name: safelyReplayable

Type: Boolean

Default: false

If true, Safely Replayable specifies that a Message is safe to send to the Remote Endpoint more than once for a single Send Action. It marks the data as safe for certain 0-RTT establishment techniques, where retransmission of the 0-RTT data may cause the remote application to receive the Message multiple times.

For protocols that do not protect against duplicated messages, e.g., UDP, all messages need to be marked as Safely Replayable. To enable protocol selection to choose such a protocol, Safely Replayable needs to be added to the TransportProperties passed to the Preconnection. If such a protocol was chosen, disabling Safely Replayable on individual messages MUST result in a SendError.

9.1.3.5. Final

Name: final

Type: Boolean

Default: false

If true, this indicates a Message is the last that the application will send on a Connection. This allows underlying protocols to indicate to the Remote Endpoint that the Connection has been effectively closed in the sending direction. For example, TCP-based Connections can send a FIN once a Message marked as Final has been completely sent, indicated by marking endOfMessage. Protocols that do not support signalling the end of a Connection in a given direction will ignore this property.

A Final Message must always be sorted to the end of a list of Messages. The Final property overrides Priority and any other property that would re-order Messages. If another Message is sent after a Message marked as Final has already been sent on a Connection, the Send Action for the new Message will cause a SendError Event.

9.1.3.6. Sending Corruption Protection Length

Name: msgChecksumLen

Type: Integer or Full Coverage

Default: Full Coverage

If this property is an Integer, it specifies the minimum length of the section of a sent Message, starting from byte 0, that the application requires to be delivered without corruption due to lower layer errors. It is used to specify options for simple integrity protection via checksums. A value of 0 means that no checksum needs to be calculated, and the enumerated value Full Coverage means that the entire Message needs to be protected by a checksum. Only Full Coverage is guaranteed, any other requests are advisory, which may result in Full Coverage being applied.

9.1.3.7. Reliable Data Transfer (Message)

Name: msgReliable

Type: Boolean

Default: the queried Boolean value of the Selection Property reliability (Section 6.2.1)

When true, this property specifies that a Message should be sent in such a way that the transport protocol ensures all data is received on the other side without corruption. Changing the Reliable Data Transfer property on Messages is only possible for Connections that were established enabling the Selection Property Configure Per-Message Reliability. When this is not the case, changing msgReliable will generate an error.

Disabling this property indicates that the Transport Services system may disable retransmissions or other reliability mechanisms for this particular Message, but such disabling is not guaranteed.

If it is not configured by the application before sending, this property's default value will be based on the Selection Property reliability of the Connection associated with the Send Action.

9.1.3.8. Message Capacity Profile Override

Name: msgCapacityProfile

Type: Enumeration

Default: inherited from the Connection Property connCapacityProfile (Section 8.1.6)

This enumerated property specifies the application's preferred tradeoffs for sending this Message; it is a per-Message override of the Capacity Profile connection property (see Section 8.1.6). If it is not configured by the application before sending, this property's default value will be based on the Connection Property connCapacityProfile of the Connection associated with the Send Action.

9.1.3.9. No Network-Layer Fragmentation

Name: noFragmentation

Type: Boolean

Default: false

This property specifies that a message should be sent and received without network-layer fragmentation, if possible. It can be used to avoid network layer fragmentation when transport segmentation is preferred.

This only takes effect when the transport uses a network layer that supports this functionality. When it does take effect, setting this property to true will cause the sender to avoid network-layer source fragmentation. When using IPv4, this will result in the Don't Fragment bit being set in the IP header.

Attempts to send a message with this property that result in a size greater than the transport's current estimate of its maximum packet size (`singularTransmissionMsgMaxLen`) can result in transport segmentation when permitted, or in a `SendError`.

Note: `noSegmentation` should be used when it is desired to only send a message within a single network packet.

9.1.3.10. No Segmentation

Name: `noSegmentation`

Type: Boolean

Default: false

When set to true, this property requests the transport layer to not provide segmentation of messages larger than the maximum size permitted by the network layer, and also to avoid network-layer source fragmentation of messages. When running over IPv4, setting this property to true can result in a sending endpoint setting the Don't Fragment bit in the IPv4 header of packets generated by the transport layer. An attempt to send a message that results in a size greater than the transport's current estimate of its maximum packet size (`singularTransmissionMsgMaxLen`) will result in a `SendError`. This only takes effect when the transport and network layer support this functionality.

9.2. Sending Data

Once a Connection has been established, it can be used for sending Messages. By default, Send enqueues a complete Message, and takes optional per-Message properties (see Section 9.2.1). All Send actions are asynchronous, and deliver Events (see Section 9.2.2). Sending partial Messages for streaming large data is also supported (see Section 9.2.3).

Messages are sent on a Connection using the Send action:

```
Connection.Send(messageData, messageContext?, endOfMessage?)
```

where messageData is the data object to send, and messageContext allows adding Message Properties, identifying Send Events related to a specific Message or inspecting meta-data related to the Message sent (see Section 9.1.1).

The optional endOfMessage parameter supports partial sending and is described in Section 9.2.3.

9.2.1. Basic Sending

The most basic form of sending on a connection involves enqueueing a single Data block as a complete Message with default Message Properties.

```
messageData := "hello"  
Connection.Send(messageData)
```

The interpretation of a Message to be sent is dependent on the implementation, and on the constraints on the Protocol Stacks implied by the Connection's transport properties. For example, a Message may be a single datagram for UDP Connections; or an HTTP Request for HTTP Connections.

Some transport protocols can deliver arbitrarily sized Messages, but other protocols constrain the maximum Message size. Applications can query the Connection Property "Maximum Message size on send" (Section 8.1.11.3) to determine the maximum size allowed for a single Message. If a Message is too large to fit in the Maximum Message Size for the Connection, the Send will fail with a SendError event (Section 9.2.2.3). For example, it is invalid to send a Message over a UDP connection that is larger than the available datagram sending size.

9.2.2. Send Events

Like all Actions in Transport Services API, the Send Action is asynchronous. There are several Events that can be delivered in response to Sending a Message. Exactly one Event (Sent, Expired, or SendError) will be delivered in response to each call to Send.

Note that if partial Sends are used (Section 9.2.3), there will still be exactly one Send Event delivered for each call to Send. For example, if a Message expired while two requests to Send data for that Message are outstanding, there will be two Expired events delivered.

The Transport Services API should allow the application to correlate which Send Action resulted in a particular Send Event. The manner in which this correlation is indicated is implementation-specific.

9.2.2.1. Sent

Connection -> Sent<messageContext>

The Sent Event occurs when a previous Send Action has completed, i.e., when the data derived from the Message has been passed down or through the underlying Protocol Stack and is no longer the responsibility of the Transport Services API. The exact disposition of the Message (i.e., whether it has actually been transmitted, moved into a buffer on the network interface, moved into a kernel buffer, and so on) when the Sent Event occurs is implementation-specific. The Sent Event contains a reference to the Message Context of the Message to which it applies.

Sent Events allow an application to obtain an understanding of the amount of buffering it creates. That is, if an application calls the Send Action multiple times without waiting for a Sent Event, it has created more buffer inside the Transport Services system than an application that always waits for the Sent Event before calling the next Send Action.

9.2.2.2. Expired

Connection -> Expired<messageContext>

The Expired Event occurs when a previous Send Action expired before completion; i.e. when the Message was not sent before its Lifetime (see Section 9.1.3.1) expired. This is separate from SendError, as it is an expected behavior for partially reliable transports. The Expired Event contains a reference to the Message Context of the Message to which it applies.

9.2.2.3. SendError

Connection -> SendError<messageContext, reason?>

A SendError occurs when a Message was not sent due to an error condition: an attempt to send a Message which is too large for the system and Protocol Stack to handle, some failure of the underlying Protocol Stack, or a set of Message Properties not consistent with the Connection's transport properties. The SendError contains a reference to the Message Context of the Message to which it applies.

9.2.3. Partial Sends

It is not always possible for an application to send all data associated with a Message in a single Send Action. The Message data may be too large for the application to hold in memory at one time, or the length of the Message may be unknown or unbounded.

Partial Message sending is supported by passing an endOfMessage boolean parameter to the Send Action. This value is always true by default, and the simpler forms of Send are equivalent to passing true for endOfMessage.

The following example sends a Message in two separate calls to Send.

```
messageContext := NewMessageContext()
messageContext.add(parameter, value)

messageData := "hel"
endOfMessage := false
Connection.Send(messageData, messageContext, endOfMessage)

messageData := "lo"
endOfMessage := true
Connection.Send(messageData, messageContext, endOfMessage)
```

All data sent with the same MessageContext object will be treated as belonging to the same Message, and will constitute an in-order series until the endOfMessage is marked.

9.2.4. Batching Sends

To reduce the overhead of sending multiple small Messages on a Connection, the application could batch several Send Actions together. This provides a hint to the system that the sending of these Messages ought to be coalesced when possible, and that sending any of the batched Messages can be delayed until the last Message in the batch is enqueued.

The semantics for starting and ending a batch can be implementation-specific, but need to allow multiple Send Actions to be enqueued.

```
Connection.StartBatch()  
Connection.Send(messageData)  
Connection.Send(messageData)  
Connection.EndBatch()
```

9.2.5. Send on Active Open: InitiateWithSend

For application-layer protocols where the Connection initiator also sends the first message, the InitiateWithSend() action combines Connection initiation with a first Message sent:

```
Connection := Preconnection.InitiateWithSend(messageData, messageContext?, timeout?)
```

Whenever possible, a messageContext should be provided to declare the Message passed to InitiateWithSend as Safely Replayable. This allows the Transport Services system to make use of 0-RTT establishment in case this is supported by the available protocol stacks. When the selected stack(s) do not support transmitting data upon connection establishment, InitiateWithSend is identical to Initiate() followed by Send().

Neither partial sends nor send batching are supported by InitiateWithSend().

The Events that may be sent after InitiateWithSend() are equivalent to those that would be sent by an invocation of Initiate() followed immediately by an invocation of Send(), with the caveat that a send failure that occurs because the Connection could not be established will not result in a SendError separate from the EstablishmentError signaling the failure of Connection establishment.

9.2.6. Priority and the Transport Services API

The Transport Services API provides two properties to allow a sender to signal the relative priority of data transmission: the Priority Message Property Section 9.1.3.2, and the Connection Priority Connection Property Section 8.1.2. These properties are designed to allow the expression and implementation of a wide variety of approaches to transmission priority in the transport and application layer, including those which do not appear on the wire (affecting only sender-side transmission scheduling) as well as those that do (e.g. [I-D.ietf-httpbis-priority]).

A Transport Services system gives no guarantees about how its expression of relative priorities will be realized. However, the Transport Services system will seek to ensure that performance of relatively-prioritized connections and messages is not worse with respect to those connections and messages than an equivalent configuration in which all prioritization properties are left at their defaults.

The Transport Services API does order Connection Priority over the Priority Message Property. In the absence of other externalities (e.g., transport-layer flow control), a priority 1 Message on a priority 0 Connection will be sent before a priority 0 Message on a priority 1 Connection in the same group.

9.3. Receiving Data

Once a Connection is established, it can be used for receiving data (unless the Direction of Communication property is set to unidirectional send). As with sending, the data is received in Messages. Receiving is an asynchronous operation, in which each call to Receive enqueues a request to receive new data from the connection. Once data has been received, or an error is encountered, an event will be delivered to complete any pending Receive requests (see Section 9.3.2). If Messages arrive at the Transport Services system before Receive requests are issued, ensuing Receive requests will first operate on these Messages before awaiting any further Messages.

9.3.1. Enqueuing Receives

Receive takes two parameters to specify the length of data that an application is willing to receive, both of which are optional and have default values if not specified.

`Connection.Receive(minIncompleteLength?, maxLength?)`

By default, Receive will try to deliver complete Messages in a single event (Section 9.3.2.1).

The application can set a `minIncompleteLength` value to indicate the smallest partial Message data size in bytes that should be delivered in response to this Receive. By default, this value is infinite, which means that only complete Messages should be delivered (see Section 9.3.2.2 and Section 9.1.2 for more information on how this is accomplished). If this value is set to some smaller value, the associated receive event will be triggered only when at least that many bytes are available, or the Message is complete with fewer bytes, or the system needs to free up memory. Applications should

always check the length of the data delivered to the receive event and not assume it will be as long as `minIncompleteLength` in the case of shorter complete Messages or memory issues.

The `maxLength` argument indicates the maximum size of a Message in bytes that the application is currently prepared to receive. The default value for `maxLength` is infinite. If an incoming Message is larger than the minimum of this size and the maximum Message size on receive for the Connection's Protocol Stack, it will be delivered via `ReceivedPartial` events (Section 9.3.2.2).

Note that `maxLength` does not guarantee that the application will receive that many bytes if they are available; the Transport Services API could return `ReceivedPartial` events with less data than `maxLength` according to implementation constraints. Note also that `maxLength` and `minIncompleteLength` are intended only to manage buffering, and are not interpreted as a receiver preference for message reordering.

9.3.2. Receive Events

Each call to `Receive` will be paired with a single Receive Event, which can be a success or an error. This allows an application to provide backpressure to the transport stack when it is temporarily not ready to receive messages.

The Transport Services API should allow the application to correlate which call to `Receive` resulted in a particular Receive Event. The manner in which this correlation is indicated is implementation-specific.

9.3.2.1. Received

Connection -> `Received<messageData, messageContext>`

A `Received` event indicates the delivery of a complete Message. It contains two objects, the received bytes as `messageData`, and the metadata and properties of the received Message as `messageContext`.

The `messageData` object provides access to the bytes that were received for this Message, along with the length of the byte array. The `messageContext` is provided to enable retrieving metadata about the message and referring to the message. The `messageContext` object is described in Section 9.1.1.

See Section 9.1.2 for handling Message framing in situations where the Protocol Stack only provides a byte-stream transport.

9.3.2.2. ReceivedPartial

Connection -> ReceivedPartial<messageData, messageContext, endOfMessage>

If a complete Message cannot be delivered in one event, one part of the Message can be delivered with a ReceivedPartial event. To continue to receive more of the same Message, the application must invoke Receive again.

Multiple invocations of ReceivedPartial deliver data for the same Message by passing the same MessageContext, until the endOfMessage flag is delivered or a ReceiveError occurs. All partial blocks of a single Message are delivered in order without gaps. This event does not support delivering discontinuous partial Messages. If, for example, Message A is divided into three pieces (A1, A2, A3) and Message B is divided into three pieces (B1, B2, B3), and preserveOrder is not Required, the ReceivedPartial may deliver them in a sequence like this: A1, B1, B2, A2, A3, B3, because the messageContext allows the application to identify the pieces as belonging to Message A and B, respectively. However, a sequence like: A1, A3 will never occur.

If the minIncompleteLength in the Receive request was set to be infinite (indicating a request to receive only complete Messages), the ReceivedPartial event may still be delivered if one of the following conditions is true:

- * the underlying Protocol Stack supports message boundary preservation, and the size of the Message is larger than the buffers available for a single message;
- * the underlying Protocol Stack does not support message boundary preservation, and the Message Framers (see Section 9.1.2) cannot determine the end of the message using the buffer space it has available; or
- * the underlying Protocol Stack does not support message boundary preservation, and no Message Framers was supplied by the application

Note that in the absence of message boundary preservation or a Message Framers, all bytes received on the Connection will be represented as one large Message of indeterminate length.

In the following example, an application only wants to receive up to 1000 bytes at a time from a Connection. If a 1500-byte message arrives, it would receive the message in two separate ReceivedPartial events.

```
Connection.Receive(1, 1000)

// Receive first 1000 bytes, message is incomplete
Connection -> ReceivedPartial<messageData(1000 bytes), messageContext, false>

Connection.Receive(1, 1000)

// Receive last 500 bytes, message is now complete
Connection -> ReceivedPartial<messageData(500 bytes), messageContext, true>
```

9.3.2.3. ReceiveError

```
Connection -> ReceiveError<messageContext, reason?>
```

A `ReceiveError` occurs when data is received by the underlying Protocol Stack that cannot be fully retrieved or parsed, and when it is useful for the application to be notified of such errors. For example, a `ReceiveError` can indicate that a `Message` (identified via the `MessageContext`) that was being partially received previously, but had not completed, encountered an error and will not be completed. This can be useful for an application, which may want to use this error as a hint to remove previously received `Message` parts from memory. As another example, if an incoming `Message` does not fulfill the Required Minimum Corruption Protection Coverage for Receiving property (see Section 8.1.1), an application can use this error as a hint to inform the peer application to adjust the Sending Corruption Protection Length property (see Section 9.1.3.6).

In contrast, internal protocol reception errors (e.g., loss causing retransmissions in TCP) are not signalled by this Event. Conditions that irrevocably lead to the termination of the Connection are signaled using `ConnectionError` (see Section 10).

9.3.3. Receive Message Properties

Each `Message Context` may contain metadata from protocols in the Protocol Stack; which metadata is available is Protocol Stack dependent. These are exposed through additional read-only `Message Properties` that can be queried from the `MessageContext` object (see Section 9.1.1) passed by the receive event. The following metadata values are supported:

9.3.3.1. UDP(-Lite)-specific Property: ECN

When available, Message metadata carries the value of the Explicit Congestion Notification (ECN) field. This information can be used for logging and debugging, and for building applications that need access to information about the transport internals for their own operation. This property is specific to UDP and UDP-Lite because these protocols do not implement congestion control, and hence expose this functionality to the application (see [RFC8293], following the guidance in [RFC8085])

9.3.3.2. Early Data

In some cases it can be valuable to know whether data was read as part of early data transfer (before connection establishment has finished). This is useful if applications need to treat early data separately, e.g., if early data has different security properties than data sent after connection establishment. In the case of TLS 1.3, client early data can be replayed maliciously (see [RFC8446]). Thus, receivers might wish to perform additional checks for early data to ensure it is safely replayable. If TLS 1.3 is available and the recipient Message was sent as part of early data, the corresponding metadata carries a flag indicating as such. If early data is enabled, applications should check this metadata field for Messages received during connection establishment and respond accordingly.

9.3.3.3. Receiving Final Messages

The Message Context can indicate whether or not this Message is the Final Message on a Connection. For any Message that is marked as Final, the application can assume that there will be no more Messages received on the Connection once the Message has been completely delivered. This corresponds to the Final property that may be marked on a sent Message, see Section 9.1.3.5.

Some transport protocols and peers do not support signaling of the Final property. Applications therefore should not rely on receiving a Message marked Final to know that the sending endpoint is done sending on a connection.

Any calls to Receive once the Final Message has been delivered will result in errors.

10. Connection Termination

A Connection can be terminated i) by the Local Endpoint (i.e., the application calls the Close, CloseGroup, Abort or AbortGroup Action), ii) by the Remote Endpoint (i.e., the remote application calls the Close, CloseGroup, Abort or AbortGroup Action), or iii) because of an error (e.g., a timeout). A local call of the Close Action will cause the Connection to either send a Closed Event or a ConnectionError Event, and a local call of the CloseGroup Action will cause all of the Connections in the group to either send a Closed Event or a ConnectionError Event. A local call of the Abort Action will cause the Connection to send a ConnectionError Event, indicating local Abort as a reason, and a local call of the AbortGroup Action will cause all of the Connections in the group to send a ConnectionError Event, indicating local Abort as a reason.

Remote Action calls map to Events similar to local calls (e.g., a remote Close causes the Connection to either send a Closed Event or a ConnectionError Event), but, different from local Action calls, it is not guaranteed that such Events will indeed be invoked. When an application needs to free resources associated with a Connection, it should therefore not rely on the invocation of such Events due to termination calls from the Remote Endpoint, but instead use the local termination Actions.

Close terminates a Connection after satisfying all the requirements that were specified regarding the delivery of Messages that the application has already given to the Transport Services system. Upon successfully satisfying all these requirements, the Connection will send the Closed Event. For example, if reliable delivery was requested for a Message handed over before calling Close, the Closed Event will signify that this Message has indeed been delivered. This Action does not affect any other Connection in the same Connection Group.

Connection.Close()

The Closed Event informs the application that a Close Action has successfully completed, or that the Remote Endpoint has closed the Connection. There is no guarantee that a remote Close will be signaled.

Connection -> Closed<>

Abort terminates a Connection without delivering any remaining Messages. This action does not affect any other Connection that is entangled with this one in a Connection Group. When the Abort Action has finished, the Connection will send a ConnectionError Event, indicating local Abort as a reason.

Connection.Abort()

CloseGroup gracefully terminates a Connection and any other Connections in the same Connection Group. For example, all of the Connections in a group might be streams of a single session for a multistreaming protocol; closing the entire group will close the underlying session. See also Section 7.4. All Connections in the group will send a Closed Event when the CloseGroup Action was successful. As with Close, any Messages remaining to be processed on a Connection will be handled prior to closing.

Connection.CloseGroup()

AbortGroup terminates a Connection and any other Connections that are in the same Connection Group without delivering any remaining Messages. When the AbortGroup Action has finished, all Connections in the group will send a ConnectionError Event, indicating local Abort as a reason.

Connection.AbortGroup()

A ConnectionError informs the application that: 1) data could not be delivered to the peer after a timeout, or 2) the Connection has been aborted (e.g., because the peer has called Abort). There is no guarantee that an Abort from the peer will be signaled.

Connection -> ConnectionError<reason?>

11. Connection State and Ordering of Operations and Events

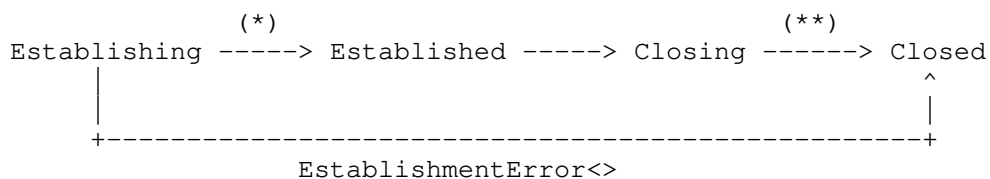
This Transport Services API is designed to be independent of an implementation's concurrency model. The details of how exactly actions are handled, and how events are dispatched, are implementation dependent.

Each transition of connection state is associated with one of more events:

- * Ready<> occurs when a Connection created with Initiate() or InitiateWithSend() transitions to Established state.

- * `ConnectionReceived<>` occurs when a Connection created with `Listen()` transitions to Established state.
- * `RendezvousDone<>` occurs when a Connection created with `Rendezvous()` transitions to Established state.
- * `Closed<>` occurs when a Connection transitions to Closed state without error.
- * `EstablishmentError<>` occurs when a Connection created with `Initiate()` transitions from Establishing state to Closed state due to an error.
- * `ConnectionError<>` occurs when a Connection transitions to Closed state due to an error in all other circumstances.

The following diagram shows the possible states of a Connection and the events that occur upon a transition from one state to another.



(*) `Ready<>`, `ConnectionReceived<>`, `RendezvousDone<>`

(**) `Closed<>`, `ConnectionError<>`

Figure 2: Connection State Diagram

The Transport Services API provides the following guarantees about the ordering of operations:

- * `Sent<>` events will occur on a Connection in the order in which the Messages were sent (i.e., delivered to the kernel or to the network interface, depending on implementation).
- * `Received<>` will never occur on a Connection before it is Established; i.e. before a `Ready<>` event on that Connection, or a `ConnectionReceived<>` or `RendezvousDone<>` containing that Connection.

- * No events will occur on a Connection after it is Closed; i.e., after a Closed<> event, an EstablishmentError<> or ConnectionError<> will not occur on that connection. To ensure this ordering, Closed<> will not occur on a Connection while other events on the Connection are still locally outstanding (i.e., known to the Transport Services API and waiting to be dealt with by the application).

12. IANA Considerations

RFC-EDITOR: Please remove this section before publication.

This document has no Actions for IANA. Later versions of this document may create IANA registries for generic transport property names and transport property namespaces (see Section 4.1).

13. Privacy and Security Considerations

This document describes a generic API for interacting with a Transport Services system. Part of this API includes configuration details for transport security protocols, as discussed in Section 6.3. It does not recommend use (or disuse) of specific algorithms or protocols. Any API-compatible transport security protocol ought to work in a Transport Services system. Security considerations for these protocols are discussed in the respective specifications.

The described API is used to exchange information between an application and the Transport Services system. While it is not necessarily expected that both systems are implemented by the same authority, it is expected that the Transport Services system implementation is either provided as a library that is selected by the application from a trusted party, or that it is part of the operating system that the application also relies on for other tasks.

In either case, the Transport Services API is an internal interface that is used to change information locally between two systems. However, as the Transport Services system is responsible for network communication, it is in the position to potentially share any information provided by the application with the network or another communication peer. Most of the information provided over the Transport Services API are useful to configure and select protocols and paths and are not necessarily privacy sensitive. Still, some information could be privacy sensitive because it might reveal usage characteristics and habits of the user of an application.

Of course any communication over a network reveals usage characteristics, as all packets, as well as their timing and size, are part of the network-visible wire image [RFC8546]. However, the selection of a protocol and its configuration also impacts which information is visible, potentially in clear text, and which other entities can access it. In most cases, information provided for protocol and path selection should not directly translate to information that can be observed by network devices on the path. However, there might be specific configuration information that is intended for path exposure, e.g., a DiffServ codepoint setting, that is either provided directly by the application or indirectly configured for a traffic profile.

Applications should be aware that communication attempts can lead to more than one connection establishment. This is the case, for example, when the Transport Services system also executes name resolution, when support mechanisms such as TURN or ICE are used to establish connectivity, if protocols or paths are raised, or if a path fails and fallback or re-establishment is supported in the Transport Services system.

Applications should also take care to not assume that all data received using the Transport Services API is always complete or well-formed. Specifically, messages that are received partially Section 9.3.2.2 could be a source of truncation attacks if applications do not distinguish between partial messages and complete messages.

The Transport Services API explicitly does not require the application to resolve names, though there is a tradeoff between early and late binding of addresses to names. Early binding allows the API implementation to reduce connection setup latency, at the cost of potentially limited scope for alternate path discovery during Connection establishment, as well as potential additional information leakage about application interest when used with a resolution method (such as DNS without TLS) which does not protect query confidentiality.

These communication activities are not different from what is used today. However, the goal of a Transport Services system is to support such mechanisms as a generic service within the transport layer. This enables applications to more dynamically benefit from innovations and new protocols in the transport, although it reduces transparency of the underlying communication actions to the application itself. The Transport Services API is designed such that protocol and path selection can be limited to a small and controlled set if required by the application for functional or security purposes. Further, A Transport Services system should provide an

interface to poll information about which protocol and path is currently in use as well as provide logging about the communication events of each connection.

14. Acknowledgements

This work has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreements No. 644334 (NEAT) and No. 688421 (MAMI).

This work has been supported by Leibniz Prize project funds of DFG - German Research Foundation: Gottfried Wilhelm Leibniz-Preis 2011 (FKZ FE 570/4-1).

This work has been supported by the UK Engineering and Physical Sciences Research Council under grant EP/R04144X/1.

This work has been supported by the Research Council of Norway under its "Toppforsk" programme through the "OCARINA" project.

Thanks to Stuart Cheshire, Josh Graessley, David Schinazi, and Eric Kinnear for their implementation and design efforts, including Happy Eyeballs, that heavily influenced this work. Thanks to Laurent Chuat and Jason Lee for initial work on the Post Sockets interface, from which this work has evolved. Thanks to Maximilian Franke for asking good questions based on implementation experience and for contributing text, e.g., on multicast.

15. References

15.1. Normative References

[I-D.ietf-taps-arch]

Pauly, T., Trammell, B., Brunstrom, A., Fairhurst, G., and C. Perkins, "An Architecture for Transport Services", Work in Progress, Internet-Draft, draft-ietf-taps-arch-12, 3 January 2022, <<https://www.ietf.org/archive/id/draft-ietf-taps-arch-12.txt>>.

[RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.

[RFC2914] Floyd, S., "Congestion Control Principles", BCP 41, RFC 2914, DOI 10.17487/RFC2914, September 2000, <<https://www.rfc-editor.org/info/rfc2914>>.

- [RFC8084] Fairhurst, G., "Network Transport Circuit Breakers", BCP 208, RFC 8084, DOI 10.17487/RFC8084, March 2017, <<https://www.rfc-editor.org/info/rfc8084>>.
- [RFC8085] Eggert, L., Fairhurst, G., and G. Shepherd, "UDP Usage Guidelines", BCP 145, RFC 8085, DOI 10.17487/RFC8085, March 2017, <<https://www.rfc-editor.org/info/rfc8085>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.
- [RFC8303] Welzl, M., Tuexen, M., and N. Khademi, "On the Usage of Transport Features Provided by IETF Transport Protocols", RFC 8303, DOI 10.17487/RFC8303, February 2018, <<https://www.rfc-editor.org/info/rfc8303>>.
- [RFC8446] Rescorla, E., "The Transport Layer Security (TLS) Protocol Version 1.3", RFC 8446, DOI 10.17487/RFC8446, August 2018, <<https://www.rfc-editor.org/info/rfc8446>>.
- [RFC8981] Gont, F., Krishnan, S., Narten, T., and R. Draves, "Temporary Address Extensions for Stateless Address Autoconfiguration in IPv6", RFC 8981, DOI 10.17487/RFC8981, February 2021, <<https://www.rfc-editor.org/info/rfc8981>>.

15.2. Informative References

- [I-D.ietf-httpbis-priority] Oku, K. and L. Pardue, "Extensible Prioritization Scheme for HTTP", Work in Progress, Internet-Draft, draft-ietf-httpbis-priority-12, 17 January 2022, <<https://www.ietf.org/archive/id/draft-ietf-httpbis-priority-12.txt>>.
- [I-D.ietf-taps-impl] Brunstrom, A., Pauly, T., Enghardt, T., Tiesel, P. S., and M. Welzl, "Implementing Interfaces to Transport Services", Work in Progress, Internet-Draft, draft-ietf-taps-impl-11, 9 January 2022, <<https://www.ietf.org/archive/id/draft-ietf-taps-impl-11.txt>>.
- [RFC2474] Nichols, K., Blake, S., Baker, F., and D. Black, "Definition of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 Headers", RFC 2474, DOI 10.17487/RFC2474, December 1998, <<https://www.rfc-editor.org/info/rfc2474>>.

- [RFC2597] Heinanen, J., Baker, F., Weiss, W., and J. Wroclawski, "Assured Forwarding PHB Group", RFC 2597, DOI 10.17487/RFC2597, June 1999, <<https://www.rfc-editor.org/info/rfc2597>>.
- [RFC3246] Davie, B., Charny, A., Bennet, J.C.R., Benson, K., Le Boudec, J.Y., Courtney, W., Davari, S., Firoiu, V., and D. Stiliadis, "An Expedited Forwarding PHB (Per-Hop Behavior)", RFC 3246, DOI 10.17487/RFC3246, March 2002, <<https://www.rfc-editor.org/info/rfc3246>>.
- [RFC3261] Rosenberg, J., Schulzrinne, H., Camarillo, G., Johnston, A., Peterson, J., Sparks, R., Handley, M., and E. Schooler, "SIP: Session Initiation Protocol", RFC 3261, DOI 10.17487/RFC3261, June 2002, <<https://www.rfc-editor.org/info/rfc3261>>.
- [RFC3376] Cain, B., Deering, S., Kouvelas, I., Fenner, B., and A. Thyagarajan, "Internet Group Management Protocol, Version 3", RFC 3376, DOI 10.17487/RFC3376, October 2002, <<https://www.rfc-editor.org/info/rfc3376>>.
- [RFC4594] Babiarz, J., Chan, K., and F. Baker, "Configuration Guidelines for DiffServ Service Classes", RFC 4594, DOI 10.17487/RFC4594, August 2006, <<https://www.rfc-editor.org/info/rfc4594>>.
- [RFC4604] Holbrook, H., Cain, B., and B. Haberman, "Using Internet Group Management Protocol Version 3 (IGMPv3) and Multicast Listener Discovery Protocol Version 2 (MLDv2) for Source-Specific Multicast", RFC 4604, DOI 10.17487/RFC4604, August 2006, <<https://www.rfc-editor.org/info/rfc4604>>.
- [RFC5245] Rosenberg, J., "Interactive Connectivity Establishment (ICE): A Protocol for Network Address Translator (NAT) Traversal for Offer/Answer Protocols", RFC 5245, DOI 10.17487/RFC5245, April 2010, <<https://www.rfc-editor.org/info/rfc5245>>.
- [RFC5280] Cooper, D., Santesson, S., Farrell, S., Boeyen, S., Housley, R., and W. Polk, "Internet X.509 Public Key Infrastructure Certificate and Certificate Revocation List (CRL) Profile", RFC 5280, DOI 10.17487/RFC5280, May 2008, <<https://www.rfc-editor.org/info/rfc5280>>.
- [RFC5482] Eggert, L. and F. Gont, "TCP User Timeout Option", RFC 5482, DOI 10.17487/RFC5482, March 2009, <<https://www.rfc-editor.org/info/rfc5482>>.

- [RFC5865] Baker, F., Polk, J., and M. Dolly, "A Differentiated Services Code Point (DSCP) for Capacity-Admitted Traffic", RFC 5865, DOI 10.17487/RFC5865, May 2010, <<https://www.rfc-editor.org/info/rfc5865>>.
- [RFC7478] Holmberg, C., Hakansson, S., and G. Eriksson, "Web Real-Time Communication Use Cases and Requirements", RFC 7478, DOI 10.17487/RFC7478, March 2015, <<https://www.rfc-editor.org/info/rfc7478>>.
- [RFC7556] Anipko, D., Ed., "Multiple Provisioning Domain Architecture", RFC 7556, DOI 10.17487/RFC7556, June 2015, <<https://www.rfc-editor.org/info/rfc7556>>.
- [RFC7657] Black, D., Ed. and P. Jones, "Differentiated Services (Diffserv) and Real-Time Communication", RFC 7657, DOI 10.17487/RFC7657, November 2015, <<https://www.rfc-editor.org/info/rfc7657>>.
- [RFC8095] Fairhurst, G., Ed., Trammell, B., Ed., and M. Kuehlewind, Ed., "Services Provided by IETF Transport Protocols and Congestion Control Mechanisms", RFC 8095, DOI 10.17487/RFC8095, March 2017, <<https://www.rfc-editor.org/info/rfc8095>>.
- [RFC8229] Pauly, T., Touati, S., and R. Mantha, "TCP Encapsulation of IKE and IPsec Packets", RFC 8229, DOI 10.17487/RFC8229, August 2017, <<https://www.rfc-editor.org/info/rfc8229>>.
- [RFC8260] Stewart, R., Tuexen, M., Loreto, S., and R. Seggelmann, "Stream Schedulers and User Message Interleaving for the Stream Control Transmission Protocol", RFC 8260, DOI 10.17487/RFC8260, November 2017, <<https://www.rfc-editor.org/info/rfc8260>>.
- [RFC8293] Ghanwani, A., Dunbar, L., McBride, M., Bannai, V., and R. Krishnan, "A Framework for Multicast in Network Virtualization over Layer 3", RFC 8293, DOI 10.17487/RFC8293, January 2018, <<https://www.rfc-editor.org/info/rfc8293>>.
- [RFC8445] Keranen, A., Holmberg, C., and J. Rosenberg, "Interactive Connectivity Establishment (ICE): A Protocol for Network Address Translator (NAT) Traversal", RFC 8445, DOI 10.17487/RFC8445, July 2018, <<https://www.rfc-editor.org/info/rfc8445>>.

- [RFC8489] Petit-Huguenin, M., Salgueiro, G., Rosenberg, J., Wing, D., Mahy, R., and P. Matthews, "Session Traversal Utilities for NAT (STUN)", RFC 8489, DOI 10.17487/RFC8489, February 2020, <<https://www.rfc-editor.org/info/rfc8489>>.
- [RFC8546] Trammell, B. and M. Kuehlewind, "The Wire Image of a Network Protocol", RFC 8546, DOI 10.17487/RFC8546, April 2019, <<https://www.rfc-editor.org/info/rfc8546>>.
- [RFC8622] Bless, R., "A Lower-Effort Per-Hop Behavior (LE PHB) for Differentiated Services", RFC 8622, DOI 10.17487/RFC8622, June 2019, <<https://www.rfc-editor.org/info/rfc8622>>.
- [RFC8656] Reddy, T., Ed., Johnston, A., Ed., Matthews, P., and J. Rosenberg, "Traversal Using Relays around NAT (TURN): Relay Extensions to Session Traversal Utilities for NAT (STUN)", RFC 8656, DOI 10.17487/RFC8656, February 2020, <<https://www.rfc-editor.org/info/rfc8656>>.
- [RFC8699] Islam, S., Welzl, M., and S. Gjessing, "Coupled Congestion Control for RTP Media", RFC 8699, DOI 10.17487/RFC8699, January 2020, <<https://www.rfc-editor.org/info/rfc8699>>.
- [RFC8838] Ivov, E., Uberti, J., and P. Saint-Andre, "Trickle ICE: Incremental Provisioning of Candidates for the Interactive Connectivity Establishment (ICE) Protocol", RFC 8838, DOI 10.17487/RFC8838, January 2021, <<https://www.rfc-editor.org/info/rfc8838>>.
- [RFC8899] Fairhurst, G., Jones, T., Tüxen, M., Rüngeler, I., and T. Völker, "Packetization Layer Path MTU Discovery for Datagram Transports", RFC 8899, DOI 10.17487/RFC8899, September 2020, <<https://www.rfc-editor.org/info/rfc8899>>.
- [RFC8922] Enghardt, T., Pauly, T., Perkins, C., Rose, K., and C. Wood, "A Survey of the Interaction between Security Protocols and Transport Services", RFC 8922, DOI 10.17487/RFC8922, October 2020, <<https://www.rfc-editor.org/info/rfc8922>>.
- [RFC8923] Welzl, M. and S. Gjessing, "A Minimal Set of Transport Services for End Systems", RFC 8923, DOI 10.17487/RFC8923, October 2020, <<https://www.rfc-editor.org/info/rfc8923>>.

[TCP-COUPLING]

Islam, S., Welzl, M., Hiorth, K., Hayes, D., Armitage, G., and S. Gjessing, "ctrlTCP: Reducing Latency through Coupled, Heterogeneous Multi-Flow TCP Congestion Control", IEEE INFOCOM Global Internet Symposium (GI) workshop (GI 2018) , 2018.

Appendix A. Implementation Mapping

The way the concepts from this abstract API map into concrete APIs in a given language on a given platform largely depends on the features and norms of the language and the platform. Actions could be implemented as functions or method calls, for instance, and Events could be implemented via event queues, handler functions or classes, communicating sequential processes, or other asynchronous calling conventions.

A.1. Types

The basic types mentioned in Section 1.1 typically have natural correspondences in practical programming languages, perhaps constrained by implementation-specific limitations. For example:

- * An Integer can typically be represented in C by an int or long, subject to the underlying platform's ranges for each.
- * In C, a Tuple may be represented as a struct with one member for each of the value types in the ordered grouping. In Python, by contrast, a Tuple may be represented natively as a tuple, a sequence of dynamically-typed elements.
- * A Collection may be represented as a std::set in C++ or as a set in Python. In C, it may be represented as an array or as a higher-level data structure with appropriate accessors defined.

The objects described in Section 1.1 can similarly be represented in different ways depending on which programming language is used. Objects like Preconnections, Connections, and Listeners can be long-lived, and benefit from using object-oriented constructs. Note that in C, these objects may need to provide a way to release or free their underlying memory when the application is done using them. For example, since a Preconnection can be used to initiate multiple Connections, it is the responsibility of the application to clean up the Preconnection memory if necessary.

A.2. Events and Errors

This specification treats Events and Errors similarly. Errors, just as any other Events, may occur asynchronously in network applications. However, implementations of this API may report Errors synchronously, according to the error handling idioms of the implementation platform, where they can be immediately detected, such as by generating an exception when attempting to initiate a connection with inconsistent Transport Properties. An error can provide an optional reason to the application with further details about why the error occurred.

A.3. Time Duration

Time duration types are implementation-specific. For instance, it could be a number of seconds, number of milliseconds, or a struct timeval in C or a user-defined Duration class in C++.

Appendix B. Convenience Functions

B.1. Adding Preference Properties

As Selection Properties of type Preference will be set on a TransportProperties object quite frequently, implementations can provide special actions for adding each preference level i.e, TransportProperties.Set(some_property, avoid) is equivalent to TransportProperties.Avoid(some_property):

```
TransportProperties.Require(property)
TransportProperties.Prefer(property)
TransportProperties.Ignore(property)
TransportProperties.Avoid(property)
TransportProperties.Prohibit(property)
```

B.2. Transport Property Profiles

To ease the use of the Transport Services API specified by this document, implementations can provide a mechanism to create Transport Property objects (see Section 6.2) that are pre-configured with frequently used sets of properties; the following are in common use in current applications:

B.2.1. reliable-inorder-stream

This profile provides reliable, in-order transport service with congestion control. TCP is an example of a protocol that provides this service. It should consist of the following properties:

Property	Value
reliability	require
preserveOrder	require
congestionControl	require
preserveMsgBoundaries	ignore

Table 2: reliable-inorder-stream preferences

B.2.2. reliable-message

This profile provides message-preserving, reliable, in-order transport service with congestion control. SCTP is an example of a protocol that provides this service. It should consist of the following properties:

Property	Value
reliability	require
preserveOrder	require
congestionControl	require
preserveMsgBoundaries	require

Table 3: reliable-message preferences

B.2.3. unreliable-datagram

This profile provides a datagram transport service without any reliability guarantee. An example of a protocol that provides this service is UDP. It consists of the following properties:

Property	Value
reliability	avoid
preserveOrder	avoid
congestionControl	ignore
preserveMsgBoundaries	require
safely replayable	true

Table 4: unreliable-datagram preferences

Applications that choose this Transport Property Profile would avoid the additional latency that could be introduced by retransmission or reordering in a transport protocol.

Applications that choose this Transport Property Profile to reduce latency should also consider setting an appropriate Capacity Profile Property, see Section 8.1.6 and might benefit from controlling checksum coverage, see Section 6.2.7 and Section 6.2.8.

Appendix C. Relationship to the Minimal Set of Transport Services for End Systems

[RFC8923] identifies a minimal set of transport services that end systems should offer. These services make all non-security-related transport features of TCP, MPTCP, UDP, UDP-Lite, SCTP and LEDBAT available that 1) require interaction with the application, and 2) do not get in the way of a possible implementation over TCP (or, with limitations, UDP). The following text explains how this minimal set is reflected in the present API. For brevity, it is based on the list in Section 4.1 of [RFC8923], updated according to the discussion in Section 5 of [RFC8923]. The present API covers all elements of this section. This list is a subset of the transport features in Appendix A of [RFC8923], which refers to the primitives in "pass 2" (Section 4) of [RFC8303] for further details on the implementation with TCP, MPTCP, UDP, UDP-Lite, SCTP and LEDBAT.

* Connect: Initiate Action (Section 7.1).

* Listen: Listen Action (Section 7.2).

- * Specify number of attempts and/or timeout for the first establishment message: timeout parameter of Initiate (Section 7.1) or InitiateWithSend Action (Section 9.2.5).
- * Disable MPTCP: multipath Property (Section 6.2.14).
- * Hand over a message to reliably transfer (possibly multiple times) before connection establishment: InitiateWithSend Action (Section 9.2.5).
- * Change timeout for aborting connection (using retransmit limit or time value): connTimeout property, using a time value (Section 8.1.3).
- * Timeout event when data could not be delivered for too long: ConnectionError Event (Section 10).
- * Suggest timeout to the peer: See "TCP-specific Properties: User Timeout Option (UTO)" (Section 8.2).
- * Notification of ICMP error message arrival: softErrorNotify (Section 6.2.17) and SoftError Event (Section 8.3.1).
- * Choose a scheduler to operate between streams of an association: connScheduler property (Section 8.1.5).
- * Configure priority or weight for a scheduler: connPriority property (Section 8.1.2).
- * "Specify checksum coverage used by the sender" and "Disable checksum when sending": msgChecksumLen property (Section 9.1.3.6) and fullChecksumSend property (Section 6.2.7).
- * "Specify minimum checksum coverage required by receiver" and "Disable checksum requirement when receiving": recvChecksumLen property (Section 8.1.1) and fullChecksumRecv property (Section 6.2.8).
- * "Specify DF field": noFragmentation property (Section 9.1.3.9).
- * Get max. transport-message size that may be sent using a non-fragmented IP packet from the configured interface: singularTransmissionMsgMaxLen property (Section 8.1.11.2).
- * Get max. transport-message size that may be received from the configured interface: recvMsgMaxLen property (Section 8.1.11.4).

- * Obtain ECN field: This is a read-only Message Property of the MessageContext object (see "UDP(-Lite)-specific Property: ECN" Section 9.3.3.1).
- * "Specify DSCP field", "Disable Nagle algorithm", "Enable and configure a Low Extra Delay Background Transfer": as suggested in Section 5.5 of [RFC8923], these transport features are collectively offered via the connCapacityProfile property (Section 8.1.6). Per-Message control ("Request not to bundle messages") is offered via the msgCapacityProfile property (Section 9.1.3.8).
- * Close after reliably delivering all remaining data, causing an event informing the application on the other side: this is offered by the Close Action with slightly changed semantics in line with the discussion in Section 5.2 of [RFC8923] (Section 10).
- * "Abort without delivering remaining data, causing an event informing the application on the other side" and "Abort without delivering remaining data, not causing an event informing the application on the other side": this is offered by the Abort action without promising that this is signaled to the other side. If it is, a ConnectionError Event will be invoked at the peer (Section 10).
- * "Reliably transfer data, with congestion control", "Reliably transfer a message, with congestion control" and "Unreliably transfer a message": data is transferred via the Send action (Section 9.2). Reliability is controlled via the reliability (Section 6.2.1) property and the msgReliable Message Property (Section 9.1.3.7). Transmitting data as a message or without delimiters is controlled via Message Framers (Section 9.1.2). The choice of congestion control is provided via the congestionControl property (Section 6.2.9).
- * Configurable Message Reliability: the msgLifetime Message Property implements a time-based way to configure message reliability (Section 9.1.3.1).
- * "Ordered message delivery (potentially slower than unordered)" and "Unordered message delivery (potentially faster than ordered)": these two transport features are controlled via the Message Property msgOrdered (Section 9.1.3.3).

- * Request not to delay the acknowledgement (SACK) of a message: should the protocol support it, this is one of the transport features the Transport Services system can apply when an application uses the connCapacityProfile Property (Section 8.1.6) or the msgCapacityProfile Message Property (Section 9.1.3.8) with value Low Latency/Interactive.
- * Receive data (with no message delimiting): Receive Action (Section 9.3) and Received Event (Section 9.3.2.1).
- * Receive a message: Receive Action (Section 9.3) and Received Event (Section 9.3.2.1), using Message Framers (Section 9.1.2).
- * Information about partial message arrival: Receive Action (Section 9.3) and ReceivedPartial Event (Section 9.3.2.2).
- * Notification of send failures: Expired Event (Section 9.2.2.2) and SendError Event (Section 9.2.2.3).
- * Notification that the stack has no more user data to send: applications can obtain this information via the Sent Event (Section 9.2.2.1).
- * Notification to a receiver that a partial message delivery has been aborted: ReceiveError Event (Section 9.3.2.3).
- * Notification of Excessive Retransmissions (early warning below abortion threshold): SoftError Event (Section 8.3.1).

Authors' Addresses

Brian Trammell (editor)
Google Switzerland GmbH
Gustav-Gull-Platz 1
CH- 8004 Zurich
Switzerland
Email: ietf@trammell.ch

Michael Welzl (editor)
University of Oslo
PO Box 1080 Blindern
0316 Oslo
Norway
Email: michawe@ifi.uio.no

Theresa Enghardt
Netflix
121 Albright Way
Los Gatos, CA 95032,
United States of America
Email: ietf@tenghardt.net

Godred Fairhurst
University of Aberdeen
Fraser Noble Building
Aberdeen, AB24 3UE
Email: gorry@erg.abdn.ac.uk
URI: <http://www.erg.abdn.ac.uk/>

Mirja Kuehlewind
Ericsson
Ericsson-Allee 1
Herzogenrath
Germany
Email: mirja.kuehlewind@ericsson.com

Colin Perkins
University of Glasgow
School of Computing Science
Glasgow G12 8QQ
United Kingdom
Email: csp@csp Perkins.org

Philipp S. Tiesel
SAP SE
Konrad-Zuse-Ring 10
14469 Potsdam
Germany
Email: philipp@tiesel.net

Tommy Pauly
Apple Inc.
One Apple Park Way
Cupertino, California 95014,
United States of America
Email: tpauly@apple.com