

TAPS Working Group
Internet-Draft
Intended status: Standards Track
Expires: 11 October 2021

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
C.A. Wood
Cloudflare
T. Pauly
Apple Inc.
K. Rose
Akamai Technologies, Inc.
9 April 2021

An Abstract Application Layer Interface to Transport Services
draft-ietf-taps-interface-12

Abstract

This document describes an abstract application programming interface, API, to the transport layer, following the Transport Services Architecture. It supports the asynchronous, atomic transmission of messages over transport protocols and network paths dynamically selected at runtime. It is intended to replace the traditional 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 11 October 2021.

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	4
1.1.	Terminology and Notation	5
1.2.	Specification of Requirements	6
2.	Overview of Interface Design	7
3.	API Summary	7
3.1.	Usage Examples	8
3.1.1.	Server Example	9
3.1.2.	Client Example	10
3.1.3.	Peer Example	11
3.2.	Transport Properties	12
3.2.1.	Transport Property Names	12
3.2.2.	Transport Property Types	13
3.3.	Scope of the Interface Definition	14
4.	Pre-Establishment Phase	15
4.1.	Specifying Endpoints	16
4.1.1.	Using Multicast Endpoints	17
4.1.2.	Endpoint Aliases	17
4.1.3.	Endpoint Examples	18
4.2.	Specifying Transport Properties	19
4.2.1.	Reliable Data Transfer (Connection)	22
4.2.2.	Preservation of Message Boundaries	22
4.2.3.	Configure Per-Message Reliability	22
4.2.4.	Preservation of Data Ordering	22
4.2.5.	Use 0-RTT Session Establishment with a Safely Replayable Message	23

4.2.6.	Multistream Connections in Group	23
4.2.7.	Full Checksum Coverage on Sending	23
4.2.8.	Full Checksum Coverage on Receiving	23
4.2.9.	Congestion control	24
4.2.10.	Keep alive	24
4.2.11.	Interface Instance or Type	24
4.2.12.	Provisioning Domain Instance or Type	25
4.2.13.	Use Temporary Local Address	26
4.2.14.	Multipath Transport	27
4.2.15.	Advertisement of Alternative Addresses	28
4.2.16.	Direction of communication	28
4.2.17.	Notification of ICMP soft error message arrival	29
4.2.18.	Initiating side is not the first to write	29
4.3.	Specifying Security Parameters and Callbacks	29
4.3.1.	Pre-Connection Parameters	30
4.3.2.	Connection Establishment Callbacks	31
5.	Establishing Connections	31
5.1.	Active Open: Initiate	31
5.2.	Passive Open: Listen	33
5.3.	Peer-to-Peer Establishment: Rendezvous	34
5.4.	Connection Groups	36
6.	Managing Connections	37
6.1.	Generic Connection Properties	39
6.1.1.	Required Minimum Corruption Protection Coverage for Receiving	39
6.1.2.	Connection Priority	39
6.1.3.	Timeout for Aborting Connection	40
6.1.4.	Timeout for keep alive packets	40
6.1.5.	Connection Group Transmission Scheduler	40
6.1.6.	Capacity Profile	41
6.1.7.	Policy for using Multipath Transports	42
6.1.8.	Bounds on Send or Receive Rate	43
6.1.9.	Group Connection Limit	43
6.1.10.	Isolate Session	44
6.1.11.	Read-only Connection Properties	44
6.2.	TCP-specific Properties: User Timeout Option (UTO)	45
6.2.1.	Advertised User Timeout	46
6.2.2.	User Timeout Enabled	46
6.2.3.	Timeout Changeable	46
6.3.	Connection Lifecycle Events	46
6.3.1.	Soft Errors	47
6.3.2.	Path change	47
7.	Data Transfer	47
7.1.	Messages and Framers	47
7.1.1.	Message Contexts	47
7.1.2.	Message Framers	48
7.1.3.	Message Properties	50
7.2.	Sending Data	56

7.2.1.	Basic Sending	56
7.2.2.	Send Events	57
7.2.3.	Partial Sends	58
7.2.4.	Batching Sends	58
7.2.5.	Send on Active Open: InitiateWithSend	59
7.2.6.	Priority in TAPS	59
7.3.	Receiving Data	60
7.3.1.	Enqueueing Receives	60
7.3.2.	Receive Events	61
7.3.3.	Receive Message Properties	63
8.	Connection Termination	64
9.	Connection State and Ordering of Operations and Events	65
10.	IANA Considerations	66
11.	Privacy and Security Considerations	66
12.	Acknowledgements	68
13.	References	68
13.1.	Normative References	69
13.2.	Informative References	69
Appendix A.	Implementation Mapping	73
A.1.	Types	73
A.2.	Events and Errors	73
A.3.	Time Duration	73
Appendix B.	Convenience Functions	74
B.1.	Adding Preference Properties	74
B.2.	Transport Property Profiles	74
B.2.1.	reliable-inorder-stream	74
B.2.2.	reliable-message	75
B.2.3.	unreliable-datagram	75
Appendix C.	Relationship to the Minimal Set of Transport Services for End Systems	76
Authors' Addresses	79

1. Introduction

This document specifies a modern abstract application programming interface (API) atop the high-level architecture for transport services defined in [I-D.ietf-taps-arch]. It supports the asynchronous, atomic transmission of messages over transport protocols and network paths dynamically selected at runtime. It is intended to replace the traditional BSD sockets API as the common interface to the transport layer, in environments where an endpoint selects from multiple interfaces and potential 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 and fallback mechanisms, which otherwise need to be separately implemented in each application.

It derives specific path and protocol selection properties and supported transport features from the analysis provided in [RFC8095], [RFC8923], and [RFC8922]. The design encourages implementations underneath the interface to dynamically choose a transport protocol depending on an application's choices rather than statically binding applications to a protocol at compile time. The Transport Services system implementations should provide 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. This specific transport stack choice is discouraged for general use, because it can reduce the portability.

1.1. Terminology and Notation

This 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 asynchronously; and Parameters associated with these Actions and Events.

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

* An Action creates an Object:

```
Object := Action()
```

* An Action creates an array of Objects:

```
[]Object := Action()
```

* An Action 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, ...>
```

Actions associated with no Object are Actions on the abstract interface itself; they are equivalent to Actions on a per-application global context.

We also make use of the following basic types:

- * **Boolean:** Instances take the value "true" or "false".
- * **Integer:** Instances take positive or negative numeric integer values, or sometimes special non-numeric (symbolic) values.
- * **Numeric:** Instances take positive or negative numeric values, or sometimes special non-numeric (symbolic) 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. The composition of types and their order depends on the property and is fixed for the property.
- * **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 Interface Design

The design of the interface 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 interface defined in this document provides:

- * Access to a variety of transport protocols, independent of the the Protocol Stacks that will be used at runtime, such that all common features of these protocol stacks are made available to the application in a transport-independent way to the degree possible, enabling applications written to a single API to make use of transport protocols in terms of the features they provide;
- * 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, allowing concurrent operations during establishment and supporting event-driven application interactions with the transport layer, in line with developments in modern platforms and programming languages;
- * 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, and for configuration of cryptographic identities and transport security parameters persistent across multiple Connections; and
- * Explicit support for multistreaming and multipath transport protocols, and the grouping of related Connections into Connection Groups through cloning of Connections, to allow applications to take full advantage of new transport protocols supporting these features.

3. API Summary

The Transport Services API is the basic common abstract application programming interface to the Transport Services Architecture defined in the TAPS Architecture [I-D.ietf-taps-arch].

An application primarily interacts with this API through two Objects: Preconnections and Connections. A Preconnection represents a set of properties and constraints on the selection and configuration of

paths and protocols to establish a Connection with a Remote Endpoint. A Connection represents a transport Protocol Stack on which data can be sent to and/or received from a Remote Endpoint (i.e., depending on the kind of transport, connections can be bi-directional or unidirectional). Connections can be created from Preconnections in three ways: by initiating the Preconnection (i.e., actively opening, as in a client), through listening on the Preconnection (i.e., passively opening, as in a server), or rendezvousing on the Preconnection (i.e. peer to peer establishment).

Once a Connection is established, data can be sent and received on it in the form of Messages. The interface supports the preservation of message boundaries both via explicit Protocol Stack support, and via application support through a Message Framer which 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.

Section 4, Section 5, Section 7.2, Section 7.3, and Section 8 describe the details of application interaction with Objects through Actions and Events in each phase of a Connection, following the phases (Pre-Establishment, Establishment, Data Transfer, and Termination) described in Section 4.1 of [I-D.ietf-taps-arch].

3.1. Usage Examples

The following usage examples illustrate how an application might use a Transport Services Interface 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 Framer, which can transform a sequence of Messages into a byte stream and vice versa (Section 7.1.2).

3.1.1. Server Example

This is an example of how an application might listen for incoming Connections using the Transport Services Interface, 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 Interface, 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()
```

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 host 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>

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

Connection -> Received<messageDataResponse, messageContext>

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

3.2. Transport Properties

Each application using the Transport Services Interface declares its preferences for how the transport service should operate using properties at each stage of the lifetime of a connection using Transport Properties, as defined in [I-D.ietf-taps-arch].

Transport Properties are divided into Selection, Connection, and Message Properties. Selection Properties (see The behavior of the selected protocol stack(s) when sending Messages is controlled by Message Properties (see Section 4.2) 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 6.1) 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. Section 7.1.3).

All Transport Properties, regardless of the phase in which they are used, are organized within a single namespace. This enables setting them as defaults at earlier stages and querying them in later stages:

- * Connection Properties can be set on Preconnections and Connections
- * Message Properties can be set on Preconnections, Connections and Messages
- * The effect of Selection Properties can be queried on Connections and Messages

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.

3.2.1. 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. These names serve two purposes:

- * Allowing different components of a TAPS 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 TAPS 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>), reformatted where necessary to conform to an implementation's naming conventions, are reserved for Protocol Specific Properties and MUST NOT be used for vendor or implementation-specific properties.

3.2.2. Transport Property Types

Transport Properties each have a type, which can be:

- * One of the basic types described in Section 1.1; or
- * Preference, which is an Enumeration with five possible values: Prohibit, Avoid, Ignore, Prefer, or Require. Each of these denotes a level of preference of a given property during protocol selection. (See Section 4.2.) The Preference type is used only on Preconnections, and only for Selection Properties.

3.3. Scope of the Interface Definition

This document defines a language- and platform-independent interface to 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 interface necessarily abstract.

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

We therefore make the following recommendations:

- * Actions, Events, and Errors in implementations of this interface 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.
- * Implementations of this interface SHOULD implement each Selection Property, Connection Property, and Message Context Property specified in this document. Each interface 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.

4. 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 has 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 4.1), the Selection Properties (see Section 4.2), any preconfigured Connection Properties (Section 6.1), and the security parameters (see Section 4.3):

```
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 host. 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 5.1) 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.

Transport Properties MUST always be specified while security parameters are OPTIONAL.

If Message Framers are used (see Section 7.1.2), they MUST be added to the Preconnection during pre-establishment.

4.1. 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) or a Service (string) that maps to a port

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

```
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 (string name)

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

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.

The Transport Services API 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 11.

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 5.3).

4.1.1. Using Multicast Endpoints

Specifying a multicast group address on a Local Endpoint will indicate to the transport 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()`. The accepted Connections 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.

4.1.2. 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. These identifiers MUST only be set for aliases.

```
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)
```

4.1.3. 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 4.2.11. 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)
```

Specify a Local Endpoint using a 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)
```

4.2. 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 for configuration of the detailed operation of the selected Protocol Stacks.

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.

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, implementations SHOULD prioritize Selection Properties that select paths over those that select protocols. Therefore, the transport system SHOULD race the path first, ignoring the protocol preference if a specific protocol does not work on the path.

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 6.

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 5.4 for more.

Section 6.1 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 this interface must provide sensible defaults for Selection Properties. The recommended 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 implementation, 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 recommended below are used, it is recommended to clearly document any differences.

4.2.1. 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.

4.2.2. Preservation of Message Boundaries

Name: preserveMsgBoundaries

Type: Preference

Default: Prefer

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

4.2.3. Configure Per-Message Reliability

Name: perMsgReliability

Type: Preference

Default: Ignore

This property specifies whether an application considers it useful to indicate its reliability requirements on a per-Message basis. This property applies to Connections and Connection Groups.

4.2.4. 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.

4.2.5. 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 7.1.3.4. Note that disabling this property has no effect for protocols that are not connection-oriented and do not protect against duplicated messages, e.g., UDP.

4.2.6. 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.

4.2.7. Full Checksum Coverage on Sending

Name: perMsgChecksumLenSend

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 7.1.3.6).

4.2.8. Full Checksum Coverage on Receiving

Name: perMsgChecksumLenRecv

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 6.1.1).

4.2.9. 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.

4.2.10. 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 6.1.4).

4.2.11. 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 all remote systems, 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 4.2.12) or another specific property.

4.2.12. 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 4.2.11), 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.

4.2.13. 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 [RFC4941]. 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.

4.2.14. 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 6.1.7 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 4.2.15 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 4.2.15 below) is set to false.

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.

4.2.15. 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 4.2.14 above); it only prevents proactive advertisement of addresses.

4.2.16. 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.

4.2.17. 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 6.3.1. 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].

4.2.18. 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].

4.3. 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 5.4).

4.3.1. Pre-Connection Parameters

Common 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, "secp256k1")
SecurityParameters.Set(ciphersuite, "TLS_ECDHE_ECDSA_WITH_CHACHA20_POLY1305_SHA25")
SecurityParameters.Set(signature-algorithm, "ed25519")
```

- * **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)
```

Representation of Security Parameters in implementations should parallel that chosen for Transport Property names as recommended in Section 3.3.

4.3.2. 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. Security handshake callbacks that may be invoked during connection establishment include:

- * Trust verification callback: Invoked when a Remote Endpoint's trust must be validated before the handshake protocol can continue.

```
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.

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

5. 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.

5.1. 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 this interface 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 7.2) 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 7.3) 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 7.2.5.

5.2. 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 this interface 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.

5.3. 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) [RFC5245].

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) [RFC5766]. 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.

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?>
```

5.4. Connection Groups

Entangled Connections can be created using the Clone Action:

```
Connection := Connection.Clone()
```

Calling Clone on a Connection yields a group of Connections: the parent Connection on which Clone was called, and a resulting cloned Connection. The connections within a group are "entangled" with each other, and become part of a Connection Group. Calling Clone on any of these Connections adds another Connection to the Connection Group, and so on. "Entangled" Connections share all Connection Properties except "Connection Priority" (see Section 6.1.2) . Like all other Properties, Connection Priority is copied to the new Connection when calling Clone(), but it is not entangled: Changing Connection Priority on one Connection does not change it on the other Connections in the same Connection Group.

The stack of Message Framers associated with a Connection are also copied to the cloned Connection when calling Clone. In other words, 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.

Changing one of the Connection Properties on one Connection in the group changes it for all others. Message Properties, however, are not entangled. For example, changing "Timeout for aborting Connection" (see Section 6.1.3) on one Connection in a group will automatically change this Connection Property for all Connections in the group in the same way. However, changing "Lifetime" (see Section 7.1.3.1) of a Message will only affect a single Message on a single Connection, entangled or not.

If the underlying protocol supports multi-streaming, it is natural to use this functionality to implement Clone. In that case, entangled Connections 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 connection, depending on the underlying Protocol Stack. When Clone() leads to multiple connections being opened instead of multi-streaming, 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-COUPPING] [RFC8699].

Attempts to clone a Connection can result in a CloneError:

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

The Connection Priority Connection Property operates on entangled Connections using the same approach as in Section 7.1.3.2: when allocating available network capacity among Connections in a Connection Group, sends on Connections with lower Priority values will be prioritized over sends on Connections with higher Priority values. Capacity will be shared among these Connections according to the Connection Group Transmission Scheduler property (Section 6.1.5). See Section 7.2.6 for more.

6. 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 6.1 below. Specific Protocol Properties are defined in a transport- and implementation-specific way, and MUST NOT be assumed to apply across different protocols. Attempts to set Specific Protocol Properties on a protocol stack not containing that specific protocol are simply ignored, and do not raise an error; however, too much reliance by an application on Specific Protocol Properties can significantly reduce the flexibility of a transport services implementation.

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 4.2), 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 5.4.

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 Section 7.1.3.5.
- * 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 7.3.3.3. The latter is only supported by certain transport protocols, e.g., by TCP as half-closed connection.
- * For Connections that are Established or Closing: Transport Properties that the application specified on the Preconnection. See Section 4.2. Selection properties of type "Preference" will be exposed as boolean values indicating whether or not the property applies to the selected transport.
- * For Connections that are Established, Closing, or Closed: Selection (Section 4.2) and Connection Properties (Section 6.1) of the actual protocols that were selected and instantiated. Selection Properties indicate whether or not the Connection has or offers a certain Selection Property. The actually instantiated

protocol stack might not match all Protocol Selection Properties that the application specified on the Preconnection. For example, a certain Protocol Selection Property that an application specified as Preferred might not actually be present in the chosen protocol stack because none of the currently available transport protocols had this feature.

- * For Connections that are Established, additional properties of the path(s) in use. These properties can be derived from the local provisioning domain [RFC7556], measurements by the Protocol Stack, or other sources.

6.1. 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.

6.1.1. Required Minimum Corruption Protection Coverage for Receiving

Name: `recvChecksumLen`

Type: Integer, with special value "Full Coverage"

Default: Full Coverage

This property specifies the minimum number of bytes in a received message that need to be covered by a checksum. A special value of 0 means that no checksum is permitted. A receiving Endpoint will not forward messages to the application that have less coverage. The application is responsible for handling any corruption within the non-protected part of the message [RFC8085].

6.1.2. Connection Priority

Name: `connPrio`

Type: Integer

Default: 100

This Property is a non-negative integer representing the relative inverse priority (i.e., a lower value reflects a higher priority) of this Connection relative to other Connections in the same Connection Group. It has no effect on Connections not part of a Connection

Group. As noted in Section 5.4, 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 7.2.6 for more details.

6.1.3. Timeout for Aborting Connection

Name: connTimeout

Type: Numeric, with special value "Disabled"

Default: Disabled

This property 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. The special value "Disabled" means that no timeout is scheduled.

6.1.4. Timeout for keep alive packets

Name: keepAliveTimeout

Type: Numeric, with special value "Disabled"

Default: Implementation-defined

A Transport Services system can request a protocol that supports sending keep alive packets Section 4.2.10. This property 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 6.1.3), or the special value "Disabled", will disable the sending of keep-alive packets.

6.1.5. 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 5.4. The set of schedulers can be taken from [RFC8260].

6.1.6. Capacity Profile

Name: connCapacityProfile

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, the 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 system 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 system 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 system 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 desired rate cannot be maintained across the Path. 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 system 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 7.1.3.8.

6.1.7. 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 4.2.14). 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.

6.1.8. Bounds on Send or Receive Rate

Name: minSendRate / minRecvRate / maxSendRate / maxRecvRate

Type: Numeric (with special value "Unlimited") / Numeric (with special value "Unlimited") / Numeric (with special value "Unlimited") / Numeric (with special value "Unlimited")

Default: Unlimited / Unlimited / Unlimited / Unlimited

This property specifies 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 special value "Unlimited" indicates that no bound is specified.

6.1.9. Group Connection Limit

Name: groupConnLimit

Type: Numeric (with special value "Unlimited")

Default: Unlimited

This property 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.

6.1.10. 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 entangled with it. Any state generated by this Connection will only be shared with entangled connections. 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.

6.1.11. Read-only Connection Properties

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

6.1.11.1. 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 7.1.3.4. It is given in Bytes.

6.1.11.2. 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 or transport layer segmentation at the sender. It exposes the Maximum Packet Size (MPS) as described in Datagram PLPMTUD [I-D.ietf-tsvwg-datagram-plpmtud].

6.1.11.3. Maximum Message Size on Send

Name: sendMsgMaxLen

Type: Integer

This property represents the maximum Message size that an application can send.

6.1.11.4. Maximum Message Size on Receive

Name: recvMsgMaxLen

Type: Integer

This numeric property represents the maximum Message size that an application can receive.

6.2. 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 implementation offers this feature, it has to expose an interface to it 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].

6.2.1. 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 6.1.3) value.

6.2.2. User Timeout Enabled

Name: tcp.userTimeout

Type: Boolean

Default: false

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

6.2.3. Timeout Changeable

Name: tcp.userTimeoutRecv

Type: Boolean

Default: true

This property controls whether the "Timeout for aborting Connection" (see Section 6.1.3) may be changed based on a UTO option received from the remote peer. This boolean becomes false when "Timeout for aborting Connection" (see Section 6.1.3) is used.

6.3. Connection Lifecycle Events

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

6.3.1. 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<>
```

6.3.2. 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, or a handover has been performed.

```
Connection -> PathChange<>
```

7. Data Transfer

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

7.1. 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.

7.1.1. Message Contexts

Using the `MessageContext` object, the application can set and retrieve meta-data of the message, including `Message Properties` (see Section 7.1.3) and framing meta-data (see Section 7.1.2.2). 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(scope?, parameter, value)  
PropertyValue := MessageContext.get(scope?, property)
```

To get or set Message Properties, the optional scope parameter is left empty. To get or set meta-data for a Framer, the application has to pass a reference to this Framer as the scope parameter.

For MessageContexts returned by send Events (see Section 7.2.2) and receive Events (see Section 7.3.2), the application can query information about the local and Remote Endpoint:

```
RemoteEndpoint := MessageContext.GetRemoteEndpoint()  
LocalEndpoint := MessageContext.GetLocalEndpoint()
```

7.1.2. 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).

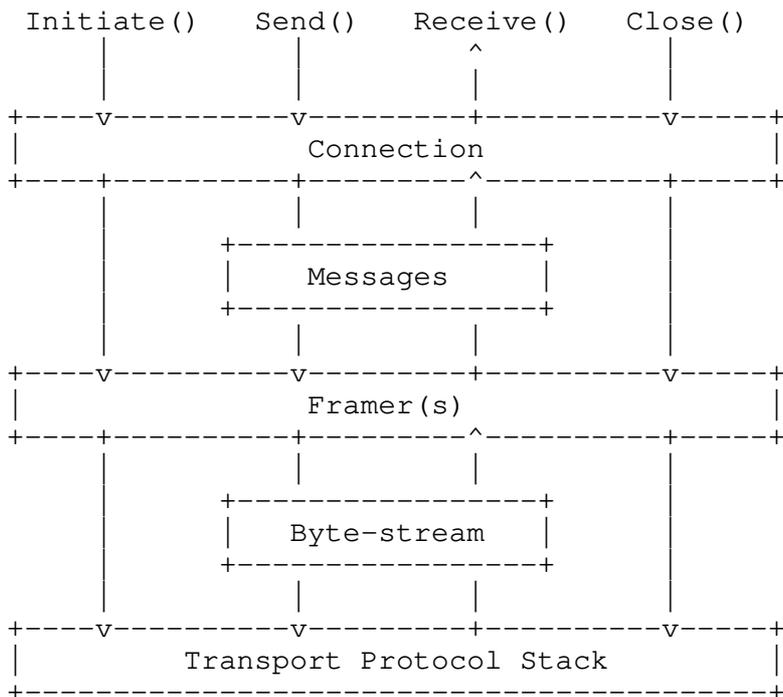


Figure 1: Protocol Stack showing a Message Framer

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].

7.1.2.1. Adding Message Framers to 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.

7.1.2.2. Framing Meta-Data

When sending Messages, applications can add Framer-specific key/value pairs to a MessageContext (Section 7.1.1). 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")
```

7.1.3. 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 7.1.1.

Message Properties are per-Message, not per-Send if partial Messages are sent (Section 7.2.3). 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, it is invalid to modify any of its properties.

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.

Connection Properties describe the default behavior for all Messages on a Connection. 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 unreliably delivered. 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:

7.1.3.1. 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 system - 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 7.1.3.7). The type and units of Lifetime are implementation-specific.

7.1.3.2. Priority

Name: msgPrio

Type: Integer (non-negative)

Default: 100

This property represents a hierarchy of priorities. It can specify 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 6.1.2. The Priority properties may interact, but can be used independently and be realized by different mechanisms; see Section 7.2.6.

7.1.3.3. Ordered

Name: msgOrdered

Type: Boolean

Default: the queried Boolean value of the Selection Property "reliability" (Section 4.2.1)

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 4.2.4, but allow out-of-order delivery for certain messages, e.g., by multiplexing independent messages onto different streams.

7.1.3.4. 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.

7.1.3.5. 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.

7.1.3.6. Sending Corruption Protection Length

Name: msgChecksumLen

Type: Integer (non-negative with special value "Full Coverage")

Default: Full Coverage

This property 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 is required, and "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.

7.1.3.7. Reliable Data Transfer (Message)

Name: msgReliable

Type: Boolean

Default: the queried Boolean value of the Selection Property "reliability" (Section 4.2.1)

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.

7.1.3.8. Message Capacity Profile Override

Name: msgCapacityProfile

Type: Enumeration

Default: inherited from the Connection Property
"connCapacityProfile" (Section 6.1.6)

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 6.1.6).

7.1.3.9. No Network-Layer Fragmentation

Name: noFragmentation

Type: Boolean

Default: false

This property specifies that a message should be sent and received as a single packet without network-layer fragmentation, if possible. 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 Don't Fragment bit to be set in the IP header, and attempts to send a message with this property set to a size greater than the transport's current estimate of its maximum packet size ("singularTransmissionMsgMaxLen") will result in a "SendError".

7.1.3.10. No Segmentation

Name: noSegmentation

Type: Boolean

Default: false

When set to true, this property requests the network layer at the sending endpoint to not fragment the packets generated by the transport layer. When running over IPv4, setting this property to true will also cause the Don't Fragment bit to be set in the IP header. When this property is set, an attempt to send a message 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.

7.2. 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 7.2.1). All Send actions are asynchronous, and deliver Events (see Section 7.2.2). Sending partial Messages for streaming large data is also supported (see Section 7.2.3).

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 7.1.1).

The optional `endOfMessage` parameter supports partial sending and is described in Section 7.2.3.

7.2.1. 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 6.1.11.3) 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 7.2.2.3). For example, it is invalid to send a Message over a UDP connection that is larger than the available datagram sending size.

7.2.2. Send Events

Like all Actions in this interface, 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 7.2.3), 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 interface 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.

7.2.2.1. 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 this interface. 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 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.

7.2.2.2. 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 7.1.3.1) 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 to which it applies.

7.2.2.3. 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` to which it applies.

7.2.3. 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.

7.2.4. 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()
```

7.2.5. 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?, timeou
```

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 InitiateError signaling the failure of Connection establishment.

7.2.6. Priority in TAPS

The Transport Services interface provides two properties to allow a sender to signal the relative priority of data transmission: the Priority Message Property Section 7.1.3.2, and the Connection Priority Connection Property Section 6.1.2. 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; for example, if a transport stack that only provides a single in-order reliable stream is selected, prioritization information can only be ignored. 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 interface 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.

7.3. 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 7.3.2). 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.

7.3.1. 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 7.3.2.1).

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 7.3.2.2 and Section 7.1.2 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 7.3.2.2).

Note that `maxLength` does not guarantee that the application will receive that many bytes if they are available; the interface 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.

7.3.2. 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 interface 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.

7.3.2.1. 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, e.g., to send replies and map responses to their requests. See Section 7.1.1 for details.

See Section 7.1.2 for handling Message framing in situations where the Protocol Stack only provides a byte-stream transport.

7.3.2.2. 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), 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 Frammer (see Section 7.1.2) 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 Frammer was supplied by the application

Note that in the absence of message boundary preservation or a Message Frammer, all bytes received on the Connection will be represented as one large Message of indeterminate length.

7.3.2.3. ReceiveError

Connection -> ReceiveError<messageContext, reason?>

A `ReceiveError` occurs when data is received by the underlying Protocol Stack that cannot be fully retrieved or parsed, or when some other indication is received that reception has failed. In contrast, conditions that irrevocably lead to the termination of the Connection are instead signaled using `ConnectionError` (see Section 8).

The `ReceiveError` event passes an optional associated `MessageContext`. This can indicate that a Message that was being partially received previously, but had not completed, encountered an error and will not be completed.

7.3.3. 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 7.1.1) passed by the receive event. The following metadata values are supported:

7.3.3.1. 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.

7.3.3.2. 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.

7.3.3.3. 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 7.1.3.5.

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.

8. Connection Termination

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. 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 that is entangled with this one in a Connection Group.

Connection.Close()

The Closed Event informs the application that the Remote Endpoint has closed the Connection. There is no guarantee that a remote Close will indeed 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.

Connection.Abort()

CloseGroup gracefully terminates a Connection and any other Connections that are entangled with this one in a 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 5.4. 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 entangled with this one in a Connection Group without delivering any remaining Messages.

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 will indeed be signaled.

Connection -> ConnectionError<reason?>

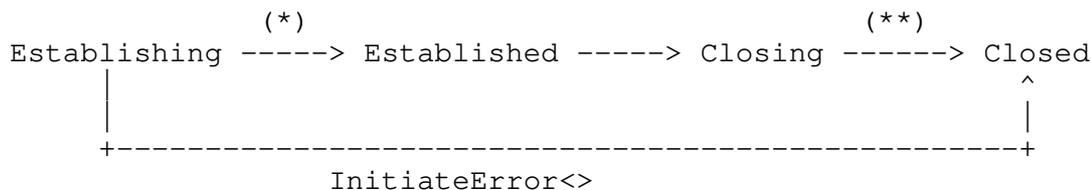
9. Connection State and Ordering of Operations and Events

This interface 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 or 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.
- * InitiateError<> 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 interface 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 InitiateError<> 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 interface and waiting to be dealt with by the application).

10. 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 3.2.1).

11. Privacy and Security Considerations

This document describes a generic API for interacting with a transport services (TAPS) system. Part of this API includes configuration details for transport security protocols, as discussed in Section 4.3. It does not recommend use (or disuse) of specific algorithms or protocols. Any API-compatible transport security protocol ought to work in a TAPS 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 7.3.2.2 could be a source of truncation attacks if applications do not distinguish between partial messages and complete messages.

The interface 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 TAPS 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, TAPS implementations 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.

12. 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.

13. References

13.1. Normative References

- [I-D.ietf-taps-arch]
Pauly, T., Trammell, B., Brunstrom, A., Fairhurst, G., Perkins, C., Tiesel, P., and C. Wood, "An Architecture for Transport Services", Work in Progress, Internet-Draft, draft-ietf-taps-arch-09, 2 November 2020, <<http://www.ietf.org/internet-drafts/draft-ietf-taps-arch-09.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>>.
- [RFC4941] Narten, T., Draves, R., and S. Krishnan, "Privacy Extensions for Stateless Address Autoconfiguration in IPv6", RFC 4941, DOI 10.17487/RFC4941, September 2007, <<https://www.rfc-editor.org/info/rfc4941>>.
- [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>>.

13.2. 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-03, 11 January 2021, <<http://www.ietf.org/internet-drafts/draft-ietf-httpbis-priority-03.txt>>.

[I-D.ietf-taps-impl]

Brunstrom, A., Pauly, T., Enghardt, T., Grinnemo, K., Jones, T., Tiesel, P., Perkins, C., and M. Welzl, "Implementing Interfaces to Transport Services", Work in Progress, Internet-Draft, draft-ietf-taps-impl-08, 2 November 2020, <<http://www.ietf.org/internet-drafts/draft-ietf-taps-impl-08.txt>>.

[I-D.ietf-tsvwg-datagram-plpmtud]

Fairhurst, G., Jones, T., Tuexen, M., Ruengeler, I., and T. Voelker, "Packetization Layer Path MTU Discovery for Datagram Transports", Work in Progress, Internet-Draft, draft-ietf-tsvwg-datagram-plpmtud-22, 10 June 2020, <<http://www.ietf.org/internet-drafts/draft-ietf-tsvwg-datagram-plpmtud-22.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>>.

- [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>>.
- [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>>.
- [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>>.
- [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/info/rfc5766>>.
- [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>>.
- [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>>.
- [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>>.
- [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 interface 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. To accommodate special values, a C function that returns a non-negative "int" on success may return -1 on failure. In Python, such a function might return "None" or raise an exception.
- * 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.

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 interface may report Errors synchronously, according to the error handling idioms of the implementation platform, where they can be immediately detected, such as by throwing 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 should 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 interface specified by this document, implementations should provide a mechanism to create Transport Property objects (see Section 4.2) that are pre-configured with frequently used sets of properties. Implementations should at least offer short-hands to specify the following property profiles:

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

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

B.2.3. unreliable-datagram

This profile provides unreliable datagram transport service. An example of a protocol that provides this service is UDP. It consists of the following properties:

Property	Value
reliability	ignore
preserveOrder	ignore
congestionControl	ignore
preserveMsgBoundaries	require
safely replayable	true

Table 4

Applications that choose this Transport Property Profile for latency reasons should also consider setting an appropriate Capacity Profile Property, see Section 6.1.6 and could benefit from controlling checksum coverage, see Section 4.2.7 and Section 4.2.8.

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 except "Notification of Excessive Retransmissions (early warning below abortion threshold)". 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 5.1).
- * Listen: "Listen" Action (Section 5.2).
- * Specify number of attempts and/or timeout for the first establishment message: "timeout" parameter of "Initiate" (Section 5.1) or "InitiateWithSend" Action (Section 7.2.5).
- * Disable MPTCP: "Multipath Transport" Property (Section 4.2.14).
- * Hand over a message to reliably transfer (possibly multiple times) before connection establishment: "InitiateWithSend" Action (Section 7.2.5).
- * Change timeout for aborting connection (using retransmit limit or time value): "Timeout for Aborting Connection" property, using a time value (Section 6.1.3).
- * Timeout event when data could not be delivered for too long: "ConnectionError" Event (Section 8).
- * Suggest timeout to the peer: "TCP-specific Properties: User Timeout Option (UTO)" (Section 6.2).
- * Notification of ICMP error message arrival: "Notification of ICMP soft error message arrival" property (Section 4.2.17).

- * Choose a scheduler to operate between streams of an association: "Connection Group Transmission Scheduler" property (Section 6.1.5).
- * Configure priority or weight for a scheduler: "Connection Priority" property (Section 6.1.2).
- * "Specify checksum coverage used by the sender" and "Disable checksum when sending": "Sending Corruption Protection Length" property (Section 7.1.3.6) and "Full Checksum Coverage on Sending" property (Section 4.2.7).
- * "Specify minimum checksum coverage required by receiver" and "Disable checksum requirement when receiving": "Required Minimum Corruption Protection Coverage for Receiving" property (Section 6.1.1) and "Full Checksum Coverage on Receiving" property (Section 4.2.8).
- * "Specify DF field": "No Network-Layer Fragmentation" property (Section 7.1.3.9).
- * "Request not to bundle messages": "No Segmentation" property (Section 7.1.3.10).
- * Get max. transport-message size that may be sent using a non-fragmented IP packet from the configured interface: "Maximum Message Size Before Fragmentation or Segmentation" property (Section 6.1.11.2).
- * Get max. transport-message size that may be received from the configured interface: "Maximum Message Size on Receive" property (Section 6.1.11.4).
- * Obtain ECN field: "UDP(-Lite)-specific Property: ECN" is a read-only Message Property of the MessageContext object (Section 7.3.3.1).
- * "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 "Capacity Profile" property (Section 6.1.6). Per-Message control is offered via the "Message Capacity Profile Override" property (Section 7.1.3.8).
- * 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 8).

- * "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 fire at the peer (Section 8).
- * "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 7.2). Reliability is controlled via the "Reliable Data Transfer (Connection)" (Section 4.2.1) property and the "Reliable Data Transfer (Message)" Message Property (Section 7.1.3.7). Transmitting data as a message or without delimiters is controlled via Message Framers (Section 7.1.2). The choice of congestion control is provided via the "Congestion control" property (Section 4.2.9).
- * Configurable Message Reliability: the "Lifetime" Message Property implements a time-based way to configure message reliability (Section 7.1.3.1).
- * "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 "Ordered" (Section 7.1.3.3).
- * 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 "Capacity Profile" Property (Section 6.1.6) or the "Message Capacity Profile Override" Message Property (Section 7.1.3.8) with value "Low Latency/Interactive".
- * Receive data (with no message delimiting): "Received" Event (Section 7.3.2.1). See Section 7.1.2 for handling Message framing in situations where the Protocol Stack only provides a byte-stream transport.
- * Receive a message: "Received" Event (Section 7.3.2.1), using Message Framers (Section 7.1.2).
- * Information about partial message arrival: "ReceivedPartial" Event (Section 7.3.2.2).
- * Notification of send failures: "Expired" Event (Section 7.2.2.2) and "SendError" Event (Section 7.2.2.3).

- * Notification that the stack has no more user data to send: applications can obtain this information via the "Sent" Event (Section 7.2.2.1).
- * Notification to a receiver that a partial message delivery has been aborted: "ReceiveError" Event (Section 7.3.2.3).

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@csperkins.org

Philipp S. Tiesel
SAP SE
Konrad-Zuse-Ring 10
14469 Potsdam
Germany

Email: philipp@tiesel.net

Christopher A. Wood
Cloudflare
101 Townsend St
San Francisco,
United States of America

Email: caw@heapingbits.net

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

Email: tpauly@apple.com

Kyle Rose
Akamai Technologies, Inc.
145 Broadway
Cambridge, MA,
United States of America

Email: krose@krose.org