

# neat

## NEAT

A New, Evolutive API and Transport-Layer Architecture for the Internet

H2020-ICT-05-2014

Project number: 644334

Deliverable D2.3

**Final Version of Core Transport System**

**Editor(s):** Naeem Khademi  
**Contributor(s):** Zdravko Bozakov, Anna Brunstrom, Øystein Dale, Dragana Damjanovic, Kristian Riktor Evensen, Gorry Fairhurst, Andreas Fischer, Karl-Johan Grinnemo, Tom Jones, Simone Mangiante, Andreas Petlund, David Ros, Irene Rüngeler, Daniel Stenberg, Michael Tüxen, Felix Weinrank and Michael Welzl

**Work Package:** 2 / Core Transport System  
**Revision:** 1.0  
**Date:** August 31, 2017  
**Deliverable type:** R (Report)  
**Dissemination level:** Public



## Abstract

This document presents the core transport system in NEAT, as used for development of the reference implementation of the NEAT System. The document describes the components necessary to realise the basic Transport Services provided by the NEAT User API, with the description of a set of NEAT building blocks and their related design choices. The design of this core transport system, which is the final product of Work Package 2, is driven by the Transport Services and API design from Task 1.4, and in close coordination with the overall NEAT architecture defined in Task 1.2.

To realise the Transport Services provided by the API, a set of transport functions has to be provided by the NEAT Core Transport System. These functions take the form of several building blocks, or *NEAT Components*, each representing an associated implementation activity. Some components are needed to ensure the basic operation of the NEAT System—e.g., a *NEAT Flow Endpoint*, a callback-based *NEAT API Framework*, the *NEAT Logic* and the functionality to *Connect to a name*. Additional components are needed for: (a) ensuring connectivity, by means of mechanisms for discovery of path support for different protocols; (b) supporting end-to-end security; (c) the ability to apply different policies to influence the decision-making process of the transport system; (d) providing other important functionalities (e.g., a user-space SCTP stack, or gathering statistics for users or system administrators).

This document updates Deliverable D2.2; in particular, the descriptions of NEAT components presented here correspond to their implementation status by the end of WP2, and as such they supersede those in D2.2.

Participant organisation name	Short name
Simula Research Laboratory AS ( <i>Coordinator</i> )	SRL
Celerway Communication AS	Celerway
EMC Information Systems International	EMC
MZ Denmark APS	Mozilla
Karlstads Universitet	KaU
Fachhochschule Münster	FHM
The University Court of the University of Aberdeen	UoA
Universitetet i Oslo	UiO
Cisco Systems France SARL	Cisco

# Contents

<b>List of Abbreviations</b>	<b>6</b>
<b>1 Introduction</b>	<b>10</b>
1.1 Overview of the NEAT Architecture	10
1.2 Overview of the services provided by the NEAT User API	12
1.3 Overview of the components required to provide the services	12
1.4 Overview of component interaction during connection setup	14
1.5 Overview of connection tear down	17
<b>2 Coding with the NEAT User API</b>	<b>17</b>
2.1 NEAT User API Tutorial	18
2.1.1 What is NEAT?	18
2.1.2 Contexts and Flows	18
2.1.3 Properties	19
2.1.4 Asynchronous API	19
2.1.5 A minimal server	19
2.1.6 A minimal client	23
2.1.7 Tying the client and server together	25
2.2 Summary of the benefits of coding with the NEAT User API	26
2.3 Python bindings for the NEAT User API	29
2.3.1 Motivation for NEAT Python bindings	29
2.3.2 Using SWIG as solution approach	30
2.3.3 Challenges and result: NEAT communication between C and Python	31
2.3.4 Summary	33
<b>3 Core Transport Functions</b>	<b>33</b>
3.1 NEAT Framework Components	33
3.1.1 NEAT Flow Endpoint	34
3.1.2 NEAT API Framework (callback)	39
3.1.3 NEAT Logic	44
3.1.4 Connect to a name	46
3.1.5 NEAT Flow Endpoint Statistics	48
3.2 NEAT Transport Components	51
3.2.1 NEAT-integrated SCTP user-space stack	51
3.2.2 Middlebox Traversal	53
3.2.3 Local flow priority	56
3.2.4 Security	59
3.3 NEAT Selection Components	61
3.3.1 Happy Eyeballs	62
3.3.2 Happy Apps (application-level feedback mechanisms)	64
3.4 NEAT Policy Components	67
3.4.1 NEAT Policy Manager	68
3.4.2 Policy Information Base (PIB)	73
3.4.3 Characteristics Information Base (CIB)	76

<b>4</b>	<b>NEAT reference material</b>	<b>80</b>
4.1	NEAT tutorial . . . . .	80
4.2	Additional online documentation . . . . .	81
4.3	Virtual machines . . . . .	81
<b>5</b>	<b>Conclusions</b>	<b>81</b>
	<b>References</b>	<b>86</b>
<b>A</b>	<b>NEAT Terminology</b>	<b>87</b>
<b>B</b>	<b>NEAT API Reference</b>	<b>89</b>
B.1	Optional arguments . . . . .	89
B.1.1	Specifying no optional arguments . . . . .	89
B.1.2	Optional argument macros . . . . .	89
B.1.3	Optional argument tags . . . . .	90
B.2	Properties . . . . .	91
B.2.1	Application property reference . . . . .	91
B.2.2	Inferred properties . . . . .	92
B.3	Callbacks . . . . .	93
B.3.1	Example callback flow . . . . .	94
B.3.2	Callback reference . . . . .	94
B.4	Error codes . . . . .	95
B.5	API functions . . . . .	96
B.5.1	neat_init_ctx . . . . .	96
B.5.2	neat_free_ctx . . . . .	96
B.5.3	neat_new_flow . . . . .	96
B.5.4	neat_set_property . . . . .	97
B.5.5	neat_get_property . . . . .	97
B.5.6	neat_open . . . . .	98
B.5.7	neat_accept . . . . .	99
B.5.8	neat_read . . . . .	100
B.5.9	neat_write . . . . .	101
B.5.10	neat_shutdown . . . . .	101
B.5.11	neat_close . . . . .	102
B.5.12	neat_abort . . . . .	102
B.5.13	neat_set_operations . . . . .	103
B.5.14	neat_change_timeout . . . . .	104
B.5.15	neat_set_primary_dest . . . . .	104
B.5.16	neat_secure_identity . . . . .	105
B.5.17	neat_set_checksum_coverage . . . . .	105
B.5.18	neat_set_qos . . . . .	106
B.5.19	neat_set_ecn . . . . .	106
B.5.20	neat_start_event_loop . . . . .	107
B.5.21	neat_stop_event_loop . . . . .	107
B.5.22	neat_get_backend_fd . . . . .	108

B.5.23	neat_get_backend_timeout . . . . .	108
B.5.24	neat_get_event_loop . . . . .	108
B.5.25	neat_get_stats . . . . .	109
B.5.26	neat_getlpaddrs . . . . .	109
B.5.27	neat_log_level . . . . .	110
B.5.28	neat_log_file . . . . .	110
<b>C</b>	<b>Main changes with respect to Deliverable D2.2</b>	<b>112</b>
<b>D</b>	<b>Paper: <i>NEAT: A Platform- and Protocol-Independent Internet Transport API</i></b>	<b>113</b>
<b>E</b>	<b>Paper: <i>On the Cost of Using Happy Eyeballs for Transport Protocol Selection</i></b>	<b>123</b>

## List of abbreviations

- AAA** Authentication, Authorisation and Accounting
- AAAA** Authentication, Authorisation, Accounting and Auditing
- API** Application Programming Interface
- BE** Best Effort
- BLEST** Blocking Estimation-based MPTCP
- CC** Congestion Control
- CCC** Coupled Congestion Controller
- CDG** CAIA Delay Gradient
- CIB** Characteristics Information Base
- CM** Congestion Manager
- DA-LBE** Deadline Aware Less than Best Effort
- DAPS** Delay-Aware Packet Scheduling
- DCCP** Datagram Congestion Control Protocol
- DNS** Domain Name System
- DNSSEC** Domain Name System Security Extensions
- DPI** Deep Packet Inspection
- DSCP** Differentiated Services Code Point
- DTLS** Datagram Transport Layer Security
- ECMP** Equal Cost Multi-Path
- EFCM** Ensemble Flow Congestion Manager
- ECN** Explicit Congestion Notification
- ENUM** Electronic Telephone Number Mapping
- E-TCP** Ensemble-TCP
- FEC** Forward Error Correction
- FLOWER** Fuzzy Lower than Best Effort
- FSE** Flow State Exchange
- FSN** Fragments Sequence Number
- GUE** Generic UDP Encapsulation
- HI** HTTP/1

**H2** HTTP/2

**HE** Happy Eyeballs

**HoLB** Head of Line Blocking

**HTTP** HyperText Transfer Protocol

**IAB** Internet Architecture Board

**ICE** Internet Connectivity Establishment

**ICMP** Internet Control Message Protocol

**IETF** Internet Engineering Task Force

**IF** Interface

**IGD-PCP** Internet Gateway Device – Port Control Protocol

**IoT** Internet of Things

**IP** Internet Protocol

**IRTF** Internet Research Task Force

**IW** Initial Window

**IW10** Initial Window of 10 segments

**JSON** JavaScript Object Notation

**KPI** Kernel Programming Interface

**LAG** Link Aggregation

**LAN** Local Area Network

**LBE** Less than Best Effort

**LEDBAT** Low Extra Delay Background Transport

**LRF** Lowest RTT First

**MBC** Model Based Control

**MID** Message Identifier

**MIF** Multiple Interfaces

**MPTCP** Multipath Transmission Control Protocol

**MPT-BM** Multipath Transport-Bufferbloat Mitigation

**MTU** Maximum Transmission Unit

**NAT** Network Address (and Port) Translation

**NEAT** New, Evolutive API and Transport-Layer Architecture

**NIC** Network Interface Card

**NUM** Network Utility Maximization

**OF** OpenFlow

**OS** Operating System

**OTIAS** Out-of-order Transmission for In-order Arrival Scheduling

**OVSDB** Open vSwitch Database

**PCP** Port Control Protocol

**PDU** Protocol Data Unit

**PHB** Per-Hop Behaviour

**PI** Policy Interface

**PIB** Policy Information Base

**PID** Proportional-Integral-Differential

**PLUS** Path Layer UDP Substrate

**PM** Policy Manager

**PMTU** Path MTU

**POSIX** Portable Operating System Interface

**PPID** Payload Protocol Identifier

**PRR** Proportional Rate Reduction

**PvD** Provisioning Domain

**QoS** Quality of Service

**QUIC** Quick UDP Internet Connections

**RACK** Recent Acknowledgement

**RFC** Request for Comments

**RTT** Round Trip Time

**RTP** Real-time Protocol

**RTSP** Real-time Streaming Protocol

**SCTP** Stream Control Transmission Protocol

**SCTP-CMT** Stream Control Transmission Protocol – Concurrent Multipath Transport

**SCTP-PF** Stream Control Transmission Protocol – Potentially Failed

**SCTP-PR** Stream Control Transmission Protocol – Partial Reliability



**SDN** Software-Defined Networking

**SDT** Secure Datagram Transport

**SIMD** Single Instruction Multiple Data

**SPUD** Session Protocol for User Datagrams

**SRTT** Smoothed RTT

**STTF** Shortest Transfer Time First

**SDP** Session Description Protocol

**SIP** Session Initiation Protocol

**SLA** Service Level Agreement

**SPUD** Session Protocol for User Datagrams

**STUN** Simple Traversal of UDP through NATs

**TCB** Transmission Control Block

**TCP** Transmission Control Protocol

**TCPINC** TCP Increased Security

**TLS** Transport Layer Security

**TSN** Transmission Sequence Number

**TTL** Time to Live

**TURN** Traversal Using Relays around NAT

**UDP** User Datagram Protocol

**UPnP** Universal Plug and Play

**URI** Uniform Resource Identifier

**VoIP** Voice over IP

**VM** Virtual Machine

**VPN** Virtual Private Network

**WAN** Wide Area Network

**WWAN** Wireless Wide Area Network

# 1 Introduction

The NEAT System aims to change the transport layer interface in a way that Internet applications may specify and select a variety of Transport Services, instead of specifying a transport protocol. The Transport Services provided by the NEAT System and the API needed to achieve this goal are outlined in Deliverable D1.3 [41]. The NEAT *Core Transport System* plays a vital role in translating those Transport Services exposed by the NEAT User API into protocol-level function calls, as well as supporting a variety of transport-layer mechanisms that provide such Transport Services<sup>1</sup>. We denote as Core Transport System the set of building blocks necessary to provide the NEAT Transport Services described in Deliverable D1.3. These building blocks include mechanisms for ensuring end-to-end connectivity, discovery of path support for protocol(s) chosen by the NEAT System, end-to-end security, ability to select different system policies depending on the application or network scenarios, and the ability to expose connection-level or system-wide statistics to an application. It also includes the support for necessary transport protocols in user-space (e.g., SCTP support by `usrctp`), and essential functionalities for using multi-streaming and multi-homing.

This document is an update to D2.2 [30] and reports the core transport system implemented in NEAT, which we term the *NEAT Library*, by the end of Work Package 2. The rest of this section provides a short overview of the NEAT architecture introduced in Deliverable D1.1 [20] (§ 1.1). It also discusses the Transport Services provided by NEAT (§ 1.2). Then, it briefly presents the building blocks, or *NEAT Components*, that comprise the NEAT core transport system (§ 1.3) and how they interact during a connection setup attempt (§ 1.4). Finally, connection closing is explained (§ 1.5).

Section 2 describes how applications can be built to use the NEAT User API. The section provides a tutorial on how to use the NEAT User API along with a simple client/server example using the NEAT library, and summarises some benefits of NEAT for application developers when compared with the traditional socket API. The section closes by presenting the Python bindings for NEAT that have been developed by the project during the last phase of WP2. Section 3 discusses each of the NEAT components in detail, identifies the Transport Services they provide as well as their dependency on other components in the NEAT System, and provides examples of operation. Section 4 provides references to related documentation and code examples. Then, Section 5 draws conclusions from this document representing the final work done in Work Package 2.

To close the document, Appendix B provides a detailed reference of the NEAT User API. Appendix C lists the main differences between D2.2 and D2.3. The paper attached in Appendix D [29] summarises some of the topics dealt with in this document, in particular the benefits offered by NEAT in different contexts, and also discusses the relation between NEAT and Internet standards activities. Finally, Appendix E contains a paper [33] with our performance evaluation results of the Happy Eyeballs mechanism used by NEAT (§ 3.3.1).

## 1.1 Overview of the NEAT Architecture

The NEAT architecture has been described in detail in Deliverable D1.1. In this section we provide a summary of such description in order to put our implementation work into context. The NEAT System is a layered architecture that provides a flexible and evolvable transport system. The applications and middleware served by the NEAT System utilise a new NEAT User API that abstracts network transport. The NEAT System can provide Transport Services in a way that allows the best transport protocol to be

<sup>1</sup>For more details about NEAT-specific terminology, please refer to Appendix A.

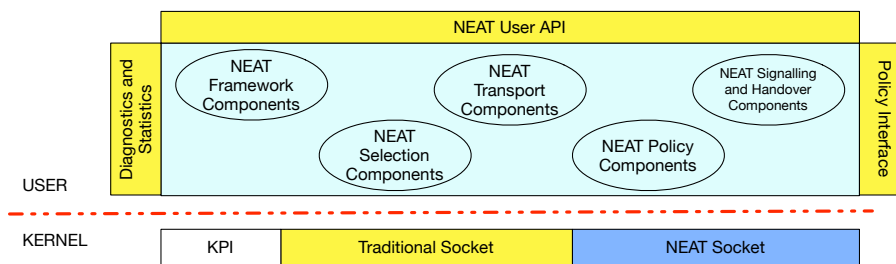


Figure 1: Component groups and interfaces used to realise the NEAT User Module.

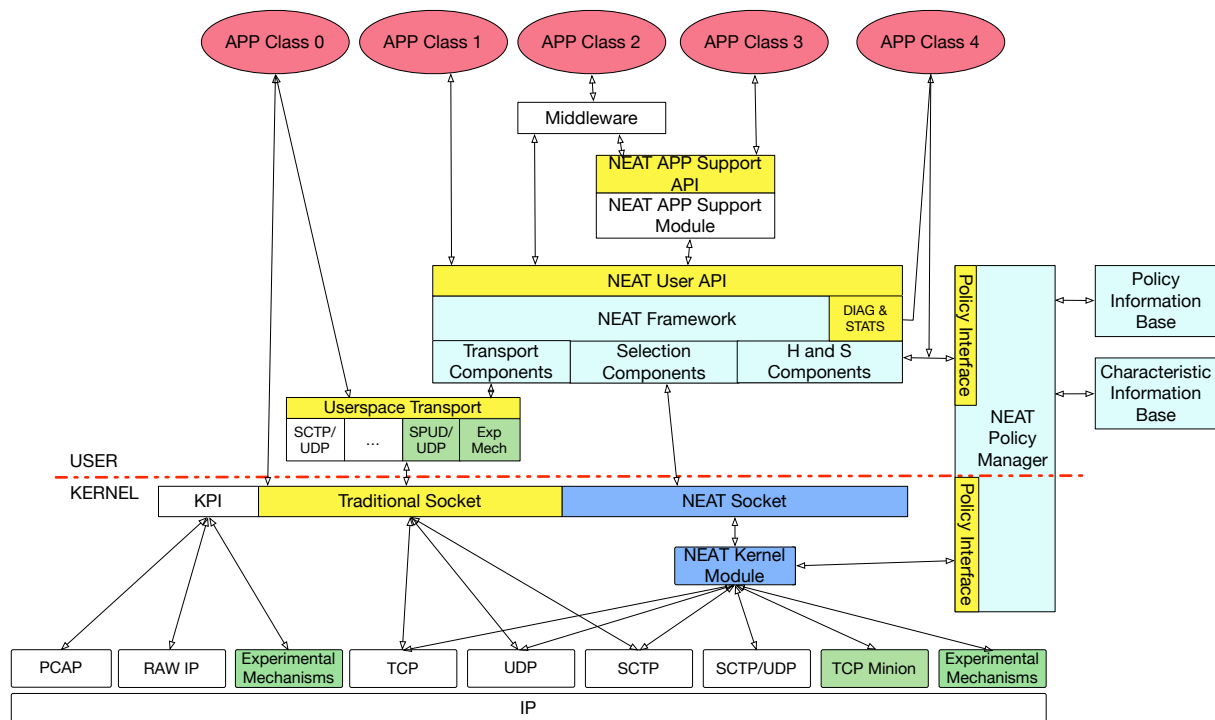


Figure 2: Components and interfaces to the NEAT System. The NEAT User Module is composed of all the blocks shown in light blue (NEAT Framework, NEAT Transport, NEAT Selection, NEAT Signalling and Handover, and Policy Components) and related APIs (NEAT User API, Policy Interface, Diagnostics and Statistics Interface).

used by an application without the application having to handle selection from application code.

The main part of the NEAT System is the NEAT User Module, depicted in Figure 1. It provides the set of components necessary to realise a Transport Service provided by the NEAT System. It is implemented in user-space and is intended to be portable across a wide range of platforms.

Figure 2 provides a more detailed overview of the different parts of the NEAT System and its interfaces. Applications access the NEAT System via a NEAT User API and its associated interfaces. The NEAT User API offers Transport Services similar to those offered by the socket API, but using an *event-driven* style of interaction. The NEAT User API provides the necessary information to allow the NEAT System to select an appropriate Transport Service.

The NEAT User API provides the interface to the NEAT User Module. This API and its associated Diagnostics and Statistics Interface are formally one part of a group of components that comprise the NEAT Framework. Other components in this group are responsible for the most basic functions of the NEAT User Module.

A group of components are responsible for the Selection of the Transport Service, these use the services of the NEAT Policy Manager, which provides high-level components that inform selection and enforce policy for decisions. The policy information is combined with the information passed via the NEAT User API and can be updated in this architecture by probing/signalling mechanisms to complete selection of the protocols and mechanisms needed to realise the required Transport Service.

The components required to configure and manage the Transport Service are also part of the NEAT User Module. Some protocols (such as TCP and UDP) are typically provided by the kernel of the platform OS. Other transport protocols are provided in user-space, but may optionally also be provided by the kernel. A key goal of the NEAT System is to offer Transport Services in the same way regardless of how the transport protocols have been implemented or how they are offered by the network stack. The NEAT User Module can utilise optional signalling components, implemented in the NEAT Signalling and Handover components.

The NEAT System can evolve to incorporate new and experimental transports. It allows applications to take advantage of new functionality as it becomes available across the Internet and will fall back and emulate features required by applications when other alternatives are not available.

The Kernel interfaces and experimental mechanisms, highlighted in Figure 2 in dark blue and green respectively, are optional components of the NEAT System.

The layered design of the NEAT System enables it to offer optimised transports to applications that would normally have to supply compatibility layers or the entire transport as a library.

## 1.2 Overview of the services provided by the NEAT User API

IETF documents from the TAPS working group [19, 39], co-authored by NEAT participants, define a Transport Service as an end-to-end service provided to an application by the transport layer, and a Transport Service Feature as a specific end-to-end feature that a Transport Service provides to its users.

Deliverable D1.2 [40] presents two sets of Transport Service Features: (a) Features derived from draft-ietf-taps-transports-usage [39]; and (b) Features derived from use cases in D1.1 [20]. Set (a) includes Transport Service Features that can be utilised using primitives and events derived from transport-protocol APIs. This includes TCP, MPTCP, SCTP, UDP and UDP-Lite protocols. Set (b) includes Transport Service Features that stem from application requirements of the use cases in D1.1 and are composed of two groups: (1) Features that require immediate action from the application; and (2) Features that require immediate action (or feedback) from NEAT.

Deliverable D1.3 [41], an update to D1.2, distilled and simplified further the abstract API of NEAT and the services it supports. The core transport system in NEAT is the minimal set of components necessary to implement such services. We summarise these components next and provide a more elaborate description in Section 3.

## 1.3 Overview of the components required to provide the services

To offer the Transport Service Features from D1.3 [41], a set of components are needed for the NEAT core transport system. Updated versions of core building blocks from D2.2 are provided in this document, presenting the final work done in WP2.

Based on the sets of NEAT Components defined in D1.1, the NEAT core components are categorised into:

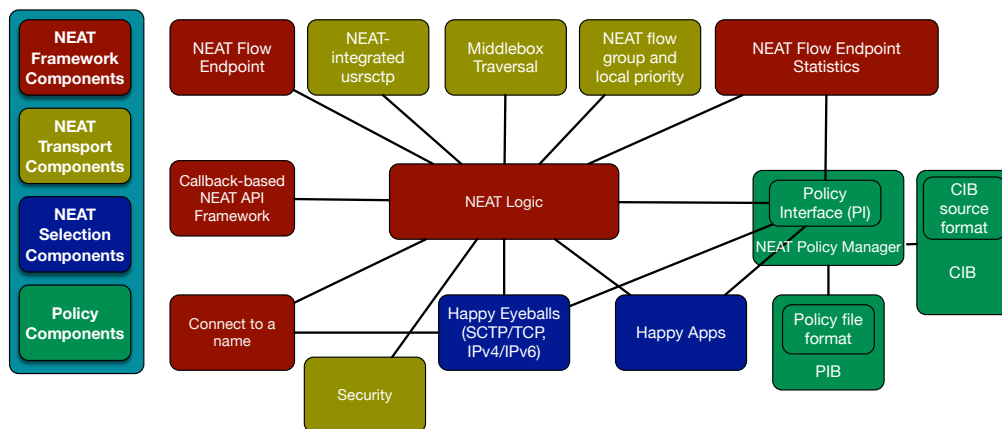


Figure 3: Building blocks of the NEAT core transport system. Each colour denotes a different component grouping.

- **NEAT Framework components:** a set of components that provide the most basic functionality required to run a NEAT System. These include the following building blocks: a *NEAT Flow Endpoint*, a callback-based *NEAT API Framework*, *NEAT Logic*, the ability to *Connect to a name* and *NEAT Flow Endpoint Statistics*.
- **NEAT Transport components:** a set of components responsible for providing the functions to configure and manage the NEAT Transport Service for a particular NEAT Flow. These include the possibility to use DTLS over UDP or SCTP as well as TLS over TCP when *Security* is being requested. In addition, the *NEAT-integrated usrsctp* component provides support for the SCTP protocol in user-land in NEAT. Achieving Peer-to-peer communication in the presence of middleboxes is achieved by integrating the WebRTC data channel into NEAT, which provides the means for *middlebox traversal*. Finally, *Local flow priority* provides a framework to use priorities among different NEAT Flows.
- **NEAT Selection components:** these components are responsible for selecting an appropriate transport endpoint and a set of protocols/mechanisms. These include building blocks for path support discovery using mechanisms at the transport layer. *Happy Eyeballs* can be done between different transport protocols (e.g., SCTP/TCP) or IP versions (IPv4/IPv6). NEAT also provides similar functionality using a feedback mechanism at the application layer named *Happy Apps*.
- **NEAT Policy components:** a set of components providing the possibility to manage and apply different policies. These building blocks include: a *NEAT Policy Manager (PM)*, a *Policy Information Base (PIB)*, a *Characteristics Information Base (CIB)* and one or multiple *CIB sources*. The NEAT Policy Manager provides a *Policy Interface (PI)* to communicate with the PIB and CIB(s) and maintains policies defined by the application developer or system developer using a predefined file format.

Figure 3 illustrates the set of core system components and their potential dependency on each other. As seen on the list above, such components are found in four out of the five component groupings in Figure 1 (Framework, Transport, Selection, and Policy); NEAT Signalling and Handover components are part of the Extended Transport System developed in Work Package 3, and therefore out of

the scope of this document<sup>2</sup>.

## 1.4 Overview of component interaction during connection setup

This section presents an overview of NEAT operations from the moment a connection setup attempt is made at the API level until the connection handle is returned to the user. This process has been depicted in Figure 4, a simplified workflow showing how NEAT components interact when opening a flow. Further details on the workings of each component will be provided in Section 3.

With regards to the complexity of the connection setup in NEAT, it is worth noting that this is no more complex than the mechanisms currently used by widely popular applications, such as state-of-the-art web browsers like Mozilla Firefox or Google Chrome, in order to establish a functioning end-to-end connection. Such mechanisms already include caching the outcome of connection attempts, trying multiple protocols (e.g., in case of QUIC or TCP Fast Open (TFO)) and so on. NEAT presents to developers of new applications a low entry barrier to implement and use advanced mechanisms similar to those deployed by such major applications.

A connection is set up in NEAT as shown in Figure 4, as follows (numbers in the arrows correspond to the step numbers below):

1. The application specifies properties of the communication with `neat_set_property` and the application calls `neat_open` in the NEAT User API.
2. The NEAT Logic sends application properties and inferred properties to the Policy Manager using a JSON object to query for feasible transport candidates, based on the destination domain name. This is done for pre-filtering purposes, to avoid running DNS lookups for non-feasible candidates since this may increase setup latency.
3. The Policy Manager finds available transport candidates based on the available policies in PIB and cached information in CIB.
4. The Policy Manager replies to NEAT Logic's query with an initial set of transport candidates eligible for name resolution (e.g., DNS lookup).
5. The NEAT Logic initiates one or more name resolution requests to the Name Resolver and the Name Resolver replies with resolved addresses, then the NEAT Logic inserts these into each candidate.
6. The NEAT Logic makes a second call to the Policy Manager, asking it to build candidates for attempting flow establishment based on the outcome of name resolution.
7. The Policy Manager builds candidates, with priorities given by policy and available CIB information. A candidate consists of the following:
  - Transport protocol
  - Interface
  - Port
  - Local address

---

<sup>2</sup>The *SDN controller integration* work, first described in Deliverable D3.1 [22], provides an example of NEAT Signalling components. See also [13].

1. Request to open flow & pass application requirements
2. Query PM about feasible transport candidates based on destination domain name
3. PM determines available transport candidates that fulfill policy (PIB) and cached information (CIB)
4. Return ranked list of feasible transport candidates as pre-filter for address resolution
5. Resolve addresses
6. Query PM about feasible transport candidates for resolved destination address
7. PM builds candidates, assigning priorities based on PIB/CIB matches
8. Return ranked list of feasible transport candidates for flow establishment
9. Do Happy Eyeballs with candidates, according to specified priorities
10. Return handle to selected transport solution
11. Cache results from Happy Eyeballs in the CIB

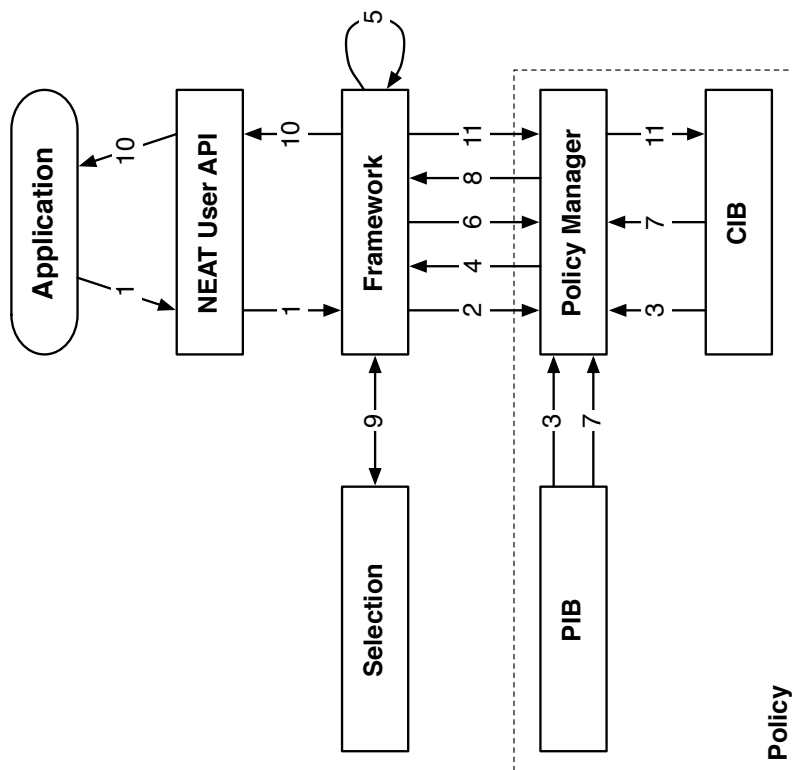


Figure 4: Simplified workflow showing how NEAT components interact when opening a flow.

- Remote address
  - Priority
  - Application properties
8. The Policy Manager returns the list of suitable candidates that the NEAT Logic should use to establish a flow. If one or more of the application properties are specified as desired (i.e., precedence 1), multiple candidates may be generated with different settings for that property.
  9. The NEAT Logic generates a list of candidates and initiates the Happy Eyeballs algorithm. The Happy Eyeballs module tries to connect each candidate. The delay between each candidate is determined from the priority of the candidate. A lower value implies a higher priority. The connection may be handled by either the operating system using its own implementation of the protocol, or using a user-space implementation of the protocol. A handle to the first connection that connects successfully and meets all required properties set by the application is returned to the NEAT Logic along with the outcome of Happy Eyeballs results for other transport candidates.
  10. The NEAT Logic starts polling the socket internally, and reports back to the application that a connection has been established using the `on_connected` callback if it has been specified using `neat_set_operations`.
  11. The NEAT Logic makes a third call to the Policy Manager, asking for the results from Happy Eyeballs to be cached in the CIB.

Now that the NEAT Flow is open for use, NEAT Logic will report that the flow is readable or writable if the respective `on_readable` or `on_writable` callbacks have been specified with `neat_set_operations`. The application can close the flow with `neat_close`.

The following provides a step-by-step example of how the mapping of application properties to the transport candidates is done by the Policy Manager internally:

- a. A NEAT application requests the `reliable_transport` and `low_latency` properties (step #1 in Figure 4). In addition, for the sake of simplicity let us assume the application has provided a `remote_IP` property, and thus no DNS resolution is necessary (i.e., steps #2–5 are not needed).
- b. The PM is asked to check the profiles and policies available in the PIB (step #6 in Figure 4).
  - i) A profile matching `reliable_transport` generates candidates for TCP and SCTP connections. The SCTP candidate is assigned a higher score than TCP.
  - ii) A low latency profile adds the socket options `TCP_NODELAY:1` and `SCTP_NODELAY:1`, respectively, to the candidates.
  - iii) The low latency policy further dictates that only wired interfaces should be used for the connection.
- c. From the CIB repository the PM infers that, for the requested destination, SCTP is only supported when going through interface `eth1` (step #7 in Figure 4).
- d. The PM returns the following ranked candidates to HE (step #8 in Figure 4):
  - i) Transport: SCTP, SCTP\_nodelay:1, interface:eth1
  - ii) Transport: TCP, TCP\_nodelay:1, interface:eth0
  - iii) Transport: TCP, TCP\_nodelay:1, interface:eth1



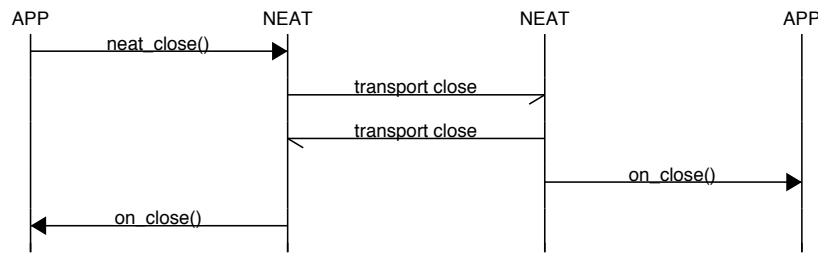


Figure 5: Schematic depiction of how NEAT end points interact when tearing down a flow. *Transport close* messages between the two end points would correspond to a FIN exchange in the case of TCP, or a SHUTDOWN / SHUTDOWN ACK / SHUTDOWN COMPLETE exchange in the case of SCTP.

## 1.5 Overview of connection tear down

This section presents an overview of NEAT operations when an established flow gets teared down. This process has been depicted in Figure 5, showing a message-sequence-diagram of involved function calls and callbacks.

The application can actively initiate the tear down of a NEAT flow by calling `neat_close` at any time. After initiating the tear down, all unsend data in the flow's outgoing buffer and all unread data in the flow's receive buffer gets discarded. Afterwards the application cannot send or receive any data from the flow by calling `neat_read` and `neat_write`, respectively. If the underlying transport protocol is connection-oriented, e.g. SCTP or TCP, NEAT will initiate the closing procedure of the particular transport protocol. After calling `neat_close`, the only callbacks that could be called from the NEAT library are `on_close`, if the tear down has been successful, and `on_error`, if an error occurred.

The successful tear down of the flow is signaled to the application via the `on_close` callback. This callback is called when the application has either actively initiated the tear down process with `neat_close` or when the remote side has closed the underlying transport connection. If the remote side has closed the transport connection, the application can still read data from the flow's receive buffer. When the `on_close` callback returns, the NEAT library frees all flow related resources and the closing procedure is finished.

## 2 Coding with the NEAT User API

Network application programmers need to understand how to program against the NEAT User API, and the way in which programs can take advantage of the services offered by a NEAT System.

Before describing the core transport functions, this section provides a tutorial that describes how an application uses the NEAT User API. It introduces the new API, which is simpler and more flexible than the existing sockets API. Key concepts are then explained, such as the role of NEAT contexts, NEAT flows and flow properties. A minimal client and server example is finally presented to illustrate how network programming tasks can be eased using this API. In the example, snippets of code illustrate the main concepts of the API. This tutorial is publicly available online<sup>3</sup> at: <http://neat.readthedocs.io/en/latest/tutorial.html>.

This section also presents an example that summarises the benefits of using the asynchronous and non-blocking NEAT User API when a developer writes an application against it.

<sup>3</sup>The version included here has been edited slightly.

Finally, bindings for the Python programming language are described. These have been developed by the NEAT project and made available on the public Github repository.

## 2.1 NEAT User API Tutorial

### 2.1.1 What is NEAT?

NEAT is a library for networked applications, intended to replace existing socket APIs with a simpler, more flexible API. Additionally, NEAT enables endpoints to make better decisions as to how to utilise the available network resources and adapts based on the current condition of the network.

With NEAT, applications are able to specify the service they want from the transport layer. NEAT will determine which of the available protocols fit the requirements of the application and tune all the relevant parameters to ensure that the application gets the desired service from the transport layer, based on knowledge about the current state of the network when this information is available.

NEAT enables applications to be written in a protocol-agnostic way, thus allowing applications to be future-proof, leveraging new protocols as they become available, with minimal to no change. Further, NEAT will try to connect with different protocols if possible, making it able to gracefully fall back to another protocol if it turns out that the most optimal protocol is unavailable, for example, because of a middlebox such as a firewall. A connection in the NEAT API will only fail if all protocols satisfying the requirements of the application are unable to connect, or if no available protocol can satisfy the requirements of the application.

Most operating systems support the same protocols. However, the same protocol may often have a slightly different API on different operating systems. NEAT provides the same API on all supported operating systems, which are currently Linux, FreeBSD, OpenBSD, NetBSD, and OS X. The availability of a protocol depends on whether the protocol is supported by the OS or if NEAT is compiled with support for a user-space stack that implements the protocol.

### 2.1.2 Contexts and Flows

The most important concept in the NEAT API is that of the *flow*. A flow is similar to a socket in the traditional Berkeley Socket API. It is a bidirectional link used to communicate between two applications, on which data may be written to or read from. Further, just like a socket, a flow uses some transport layer protocol to communicate.

However, one important difference is that a flow is not as strictly tied to the underlying transport protocol in the same way a socket is. In fact, a flow may be created without even specifying which transport protocol to use. This is not possible with a socket.

The same applies to modifying options on sockets. Setting the same kind of option on two sockets with different protocols in the traditional socket API requires `setsockopt` calls with different protocol IDs, option names, and sometimes even values with different units. The `setsockopt` calls also vary depending on what system you are on. This is not the case with NEAT. As long as the desired option is available for the protocol in use, the API for setting that option is the same for all protocols, and on all operating systems supported by NEAT.

A *context* is a common environment for multiple flows. Along with flows, it contains several services that are used by the flows internally in NEAT, such as a DNS resolver and a Happy Eyeballs implementation. Flows within a context are polled together. A flow may only belong to the context in

which it is created, and it cannot be transferred to a different context. Most applications need only one context.

### 2.1.3 Properties

Different types of applications have different requirements and desires to the services provided by the transport layer. An application for real-time communication may require the communication to have properties such as low latency, high bandwidth, quality of service, and have less strict requirements with regards to reliable delivery. Losing a packet or bit errors may be less critical to these applications. A web browser, on the other hand, might require communication that is (partially) ordered and error-free. A BitTorrent application might only require the ability to send packets to some destination with a minimum amount of effort, and not at the expense of other applications with stricter capacity requirements.

With the traditional socket API, the application requirements dictate the choice of which protocol to use. With NEAT, this is not the case. NEAT enables applications to specify the properties of the communication instead of specifying which protocol to use. Some properties may be required; other properties may be desired to support a specific application or user preference, but are not mandatory.

Based on the properties, NEAT will determine which protocols can support the requirements of the application and the options to set for each protocol. This mapping of application properties to a list of potential transport candidates is done with the help of the Policy Manager (PM) (§ 3.4.1).

Consequently, the NEAT System will try to establish a connection by trying each of the transport candidates until one connection succeeds, a method known as *Happy Eyeballing*.

The ability to specify properties instead of protocols allows applications to take advantage of available protocols where possible. By Happy Eyeballing, NEAT ensures that applications are able to cope with different network configurations, and gracefully fall back to another protocol if necessary should the most desirable protocol not be available for whatever reason.

### 2.1.4 Asynchronous API

The NEAT API is asynchronous and non-blocking. Once the execution is transferred to NEAT, it will poll the sockets internally, and, when an event happens, execute the appropriate callback in the application. This creates a more natural way of programming communicating applications than with the traditional socket API.

The three most important callbacks in the NEAT API are `on_connected`, `on_readable` and `on_writable`, which may be set per flow. The `on_connected` callback will be executed once the flow has connected to a remote endpoint, or a flow has connected to a server listening for incoming connections. The `on_writable` and `on_readable` callbacks are executed once data may be written to or read from the flow without blocking.

Figure 6 depicts a sufficient callback flow for most applications.

### 2.1.5 A minimal server

To get started using the NEAT API, we will write a small server that will send `Hello, this is NEAT!` to any client that connects to it. Later, we will write a similar client, before modifying this server so that it works with the client.

We summarise the functionality as follows:

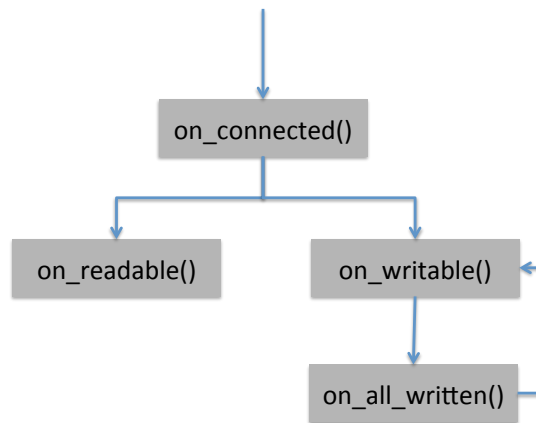


Figure 6: NEAT callback flow.

- When a client connects, start writing when the flow is writable.
- When a flow is writable, write `Hello, this is NEAT!` to it.
- When the flow has finished writing, close it.

Pay close attention to how easily this natural description can be implemented using the NEAT API.

Here are the includes that should be put on top of the file:

```
#include <neat.h>

#include <stdlib.h>
#include <stdio.h>
#include <string.h>
#include <unistd.h>
```

We will start writing the main function of our server. The first thing is to declare a few variables:

```
main(int argc, char *argv[])
{
    struct neat_ctx *ctx;
    struct neat_flow *flow;
    struct neat_flow_operations ops;
```

And initialise them:

```
ctx = neat_init_ctx();
flow = neat_new_flow(ctx);
memset(&ops, 0, sizeof(ops));
```

We are already familiar with the flow and the context. `neat_init_ctx` is used to initialise the context, and `neat_new_flow` creates a new flow within the context. The `neat_flow_operations` struct is used to tell NEAT what to do when certain events occur. We use that next to tell which function we want NEAT to call when a client connects:

```
ops.on_connected = on_connected;
neat_set_operations(ctx, flow, &ops);
```

The function `on_connected` has not been defined yet, we will do that later. Now that we have told NEAT what to do with a connecting client, we are ready to accept incoming connections.

```
if (neat_accept(ctx, flow, 5000, NULL, 0)) {  
    fprintf(stderr, "neat_accept failed\n");  
    return EXIT_FAILURE;  
}
```

This will instruct NEAT to start listening to incoming connections on port 5000. The flow passed to `neat_accept` is cloned for each accepted connection. The last two parameters are used for optional arguments. This example does not use them.

The last function call we will do in main will be the one that starts the show:

```
neat_start_event_loop(ctx, NEAT_RUN_DEFAULT);  
  
return EXIT_SUCCESS;  
}
```

When this function is called, NEAT will start doing work behind the scenes. When called with the `NEAT_RUN_DEFAULT` parameter, this function will not return until all flows have closed and all events have been handled. It is also possible to run NEAT without having NEAT capture the main loop. The final main function of our server application looks like this:

```
main(int argc, char *argv[])  
{  
    struct neat_ctx *ctx;  
    struct neat_flow *flow;  
    struct neat_flow_operations ops;  
  
    ctx = neat_init_ctx();  
    flow = neat_new_flow(ctx);  
    memset(&ops, 0, sizeof(ops));  
  
    // ops.on_readable = on_readable;  
    ops.on_connected = on_connected;  
    neat_set_operations(ctx, flow, &ops);  
  
    if (neat_accept(ctx, flow, 5000, NULL, 0)) {  
        fprintf(stderr, "neat_accept failed\n");  
        return EXIT_FAILURE;  
    }  
  
    neat_start_event_loop(ctx, NEAT_RUN_DEFAULT);  
  
    return EXIT_SUCCESS;  
}
```

It is time to start working on the callbacks that NEAT will use. The first callback we need is `on_connected`.

```
static neat_error_code  
on_connected(struct neat_flow_operations *ops)  
{
```

From the functional description above, we know that we need to write to connecting clients when this becomes possible. The callback contains a parameter that is a pointer to a `neat_flow_operations` struct, which we can use to update the active callbacks of the flow. We set the `on_writable` callback so that we can start writing when the flow becomes writable:

```
ops->on_writable = on_writable;
```

It is also good practice to set the `on_all_written` callback when setting the `on_writable` callback:

```
ops->on_all_written = on_all_written;
```

The change is applied by calling `neat_set_operations`, as in the main function:

```
neat_set_operations(ops->ctx, ops->flow, ops);

return NEAT_OK;
}
```

Next, we write the `on_writable` callback:

```
static neat_error_code
on_writable(struct neat_flow_operations *ops)
{
```

Here, we call the function that will send our message:

```
neat_write(ops->ctx, ops->flow, "Hello, this is NEAT!", 20, NULL, 0);
opCB->on_writable = NULL;
return NEAT_OK;
}
```

Here we specify the data to send and the length of the data. As with the `neat_accept` function, `neat_write` takes optional parameters. We do not need to set any optional parameters for this call either, so again we pass `NULL` and `0`.

The final callback we need to implement is the `on_all_written` callback:

```
static neat_error_code
on_all_written(struct neat_flow_operations *ops)
{
```

Here, we call `neat_close` to close the flow:

```
neat_close(ops->ctx, ops->flow);
return NEAT_OK;
}
```

This is the final piece of our server. You may now compile and run the server. You can use the tool `socat` to test it. The following output should be observed:

```
$ socat STDIO TCP:localhost:5000
Hello, this is NEAT!
$ socat STDIO SCTP:localhost:5000
Hello, this is NEAT!
```

### 2.1.6 A minimal client

Next, we implement a client that will send the message "Hi!" after connecting to a server, and then receive a reply from the server. A fair amount of the code will be similar to the server we wrote above, so you may make a copy of the code for the server and use that as a starting point for the client.

We will make two additions and one change to the main function. First, since we are connecting to a server, we change the `neat_accept` call to `neat_open` instead:

```
if (neat_open(ctx, flow, "127.0.0.1", 5000, NULL, 0)) {
    fprintf(stderr, "neat_open failed\n");
    return EXIT_FAILURE;
}
```

Next, we will specify a few properties for the flow:

```
static char *properties = "{\n\
  \"transport\": [\n\
    {\n\
      \"value\": \"SCTP\",\n\
      \"precedence\": 1\n\
    },\n\
    {\n\
      \"value\": \"TCP\",\n\
      \"precedence\": 1\n\
    }\n\
  ]\n\
}";
```

These properties will tell NEAT that it can select either SCTP or TCP as the transport protocol. The properties are applied with `neat_set_properties`, which may be done at any point between `neat_new_flow` and `neat_open`.

Finally, we add `neat_free_ctx` after `neat_start_event_loop`, so that NEAT may free any allocated resources and exit gracefully. The complete main function of the client will look like this:

```
int
main(int argc, char *argv[])
{
    struct neat_ctx *ctx;
    struct neat_flow *flow;
    struct neat_flow_operations ops;

    ctx = neat_init_ctx();
    flow = neat_new_flow(ctx);
    memset(&ops, 0, sizeof(ops));

    ops.on_connected = on_connected;
    neat_set_operations(ctx, flow, &ops);

    neat_set_property(ctx, flow, properties);

    if (neat_open(ctx, flow, "127.0.0.1", 5000, NULL, 0)) {
```

```
        fprintf(stderr, "neat_open failed\n");
        return EXIT_FAILURE;
    }

    neat_start_event_loop(ctx, NEAT_RUN_DEFAULT);

    neat_free_ctx(ctx);

    return EXIT_SUCCESS;
}
```

Leave the `on_connected` callback similar to the server.

We change the `on_writable` callback to send "Hi!" instead:

```
static neat_error_code
on_writable(struct neat_flow_operations *ops)
{
    neat_write(ops->ctx, ops->flow, "Hi!", 3, NULL, 0);
    return NEAT_OK;
}
```

The `on_all_written` callback should not close the flow, but instead stop writing and set the `on_readable` callback:

```
static neat_error_code
on_all_written(struct neat_flow_operations *ops)
{
    ops->on_readable = on_readable;
    ops->on_writable = NULL;
    neat_set_operations(ops->ctx, ops->flow, ops);
    return NEAT_OK;
}
```

Finally, we will write an `on_readable` callback for the client. We allocate some space on the stack to store the received data, and use a variable to store the length of the received message. If the `neat_read` call returns successfully, we print the message. Finally, we stop the internal event loop in NEAT, which will eventually cause the call to `neat_start_event_loop` in the main function to return. The `on_readable` callback should look like this:

```
static neat_error_code
on_readable(struct neat_flow_operations *ops)
{
    uint32_t bytes_read = 0;
    char buffer[32];

    if (neat_read(ops->ctx, ops->flow, buffer, 31, &bytes_read, NULL, 0) == NEAT_OK)
    {
        buffer[bytes_read] = 0;
        fprintf(stdout, "Read %u bytes:\n%s", bytes_read, buffer);
    }

    neat_close(ops->ctx, ops->flow);
}
```



```
    return NEAT_OK;
}
```

And there we have our finished client! You can test it with `socat`:

```
$ socat TCP-LISTEN:5000 STDIO
```

When you run the client, you should see `Hi!` show up in the output from `socat`. You can type a short message followed by pressing `return`, and it should show up in the output on the client.

### 2.1.7 Tying the client and server together

A few small changes are required on the server to make the client and server work together. In the `on_connected` callback, the server should set the `on_readable` callback instead of the `on_writable` callback. An `on_readable` callback should be added and read the incoming message from the client, and set the `on_writable` callback.

The callbacks for the updated server are:

```
static neat_error_code
on_readable(struct neat_flow_operations *ops)
{
    uint32_t bytes_read = 0;
    char buffer[32];

    if (neat_read(ops->ctx, ops->flow, buffer, 31, &bytes_read, NULL, 0) == NEAT_OK)
    {
        buffer[bytes_read] = 0;
        fprintf(stdout, "Read %u bytes:\n%s\n", bytes_read, buffer);
    }

    ops->on_readable = NULL;
    ops->on_writable = on_writable;
    ops->on_all_written = on_all_written;
    neat_set_operations(ops->ctx, ops->flow, ops);

    return NEAT_OK;
}

static neat_error_code
on_writable(struct neat_flow_operations *ops)
{
    neat_write(ops->ctx, ops->flow, "Hello, this is NEAT!", 20, NULL, 0);

    ops->on_writable = NULL;
    ops->on_readable = NULL;
    ops->on_all_written = on_all_written;
    neat_set_operations(ops->ctx, ops->flow, ops);

    return NEAT_OK;
}
```

```
static neat_error_code
on_all_written(struct neat_flow_operations *ops)
{
    ops->on_readable = NULL;
    ops->on_writable = NULL;
    ops->on_all_written = NULL;
    neat_set_operations(ops->ctx, ops->flow, ops);

    neat_close(ops->ctx, ops->flow);

    return NEAT_OK;
}

static neat_error_code
on_connected(struct neat_flow_operations *ops)
{
    ops->on_readable = on_readable;
    neat_set_operations(ops->ctx, ops->flow, ops);

    return NEAT_OK;
}
```

## 2.2 Summary of the benefits of coding with the NEAT User API

This section summarises the simplifications possible when a developer writes an application using the asynchronous and non-blocking NEAT User API.

A key simplification is that the NEAT User API offers a uniform way to access networking functionality, independent from the underlying network protocol or operating system, resulting in portable network applications. Many common network programming tasks like address resolution, buffer management, encryption, connection establishment and handling are built into the NEAT Library and can be used by any application that uses NEAT.

The callback interface is implemented using the libuv [3] library which provides asynchronous I/O across multiple platforms. Coding is also simplified because the NEAT User API executes callbacks in the application when an event from the underlying transport happens, creating a more natural and less error-prone way of network programming than with the traditional socket API. This differs from the code written against the socket API since the callbacks and event loop are handled internally, whereas the socket-based code needs to maintain the event loop and perform the polling.

Our experience with various applications ported to use NEAT shows a reduction of the networking code size by  $\approx 20\%$  for each application<sup>4</sup>, as the library streamlines a number of connection establishment steps. For example, the single function call `neat_open` requests name resolution and all other functions required before communication can start, hiding complex boilerplate code. Ported applications remain fully interoperable with regular TCP/IP-based implementations, while being able to take advantage of NEAT functions. Besides, they can benefit from support for alternative transports, when available, relieving programmers from dealing with fallbacks between protocols.

<sup>4</sup>The code for several of these applications is available at <https://github.com/NEAT-project/neat/tree/master/examples>.

To illustrate some benefits of developing applications against the NEAT API, we have developed a simple example application called FileStreamer [1]. This application sends multiple files from a server to a client using the SCTP protocol. Leveraging the multi-streaming functionality in SCTP, different files can be transferred on different SCTP streams simultaneously.

Our code contains two examples on the server side, one with NEAT (`server_neat.c`) and one without NEAT, purely written with the traditional socket API (`server_sockets.c`). The client side (file receiver) is identical in both cases—NEAT is able to operate one-sided regardless of the support for it on the other end point, which may ease deployment and simplify migration of applications.

Listings 1 and 2 provide code snippets from FileStreamer, showing how a listening socket is set up with NEAT and with the socket API, respectively.

As seen in these listings, address resolution, creating and binding to a socket and their related error handling are all done internally by NEAT. Also, using the traditional socket API, SCTP has to be selected explicitly, and the `SCTP_INITMSG` socket option has to be initialised with the number of streams the application wishes to use (lines 42 to 48 of Listing 2). A similar setting without any protocol-related parameter configuration is achieved by setting the protocol-agnostic `NEAT_TAG_STREAM_COUNT` value and the property that expresses the requirement to use SCTP using `neat_set_property` (lines 24 and 30 of Listing 1).

In addition, the socket-based code differs from the code with NEAT in that the former needs to handle itself the polling and event loop (`do_poll` function in `server_sockets.c`, not shown here).

```
1 int
2 setup_neat(const char* port, int file_count, char *file_names[])
3 {
4     int rc = 0;
5     struct neat_ctx *ctx = NULL;
6     struct neat_flow *flow = NULL;
7     struct neat_flow_operations ops;
8     NEAT_OPTARGS_DECLARE(1);
9
10    NEAT_OPTARGS_INIT();
11    memset(&ops, 0, sizeof(ops));
12
13    if ((ctx = neat_init_ctx()) == NULL)
14        return -1;
15
16    if ((flow = neat_new_flow(ctx)) == NULL) {
17        rc = -1;
18        goto error;
19    }
20
21    no_of_files = file_count;
22    files = file_names;
23
24    NEAT_OPTARG_INT(NEAT_TAG_STREAM_COUNT, file_count);
25
26    ops.on_connected = on_connected;
27    ops.on_error = on_error;
28    neat_set_operations(ctx, flow, &ops);
```

```
29
30     neat_set_property(ctx, flow, "{\"transport\": [{\"value\": \"SCTP\", \"
      precedence\": 2}]}");
31
32     if (neat_accept(ctx, flow, 5001, NEAT_OPTARGS, NEAT_OPTARGS_COUNT) != NEAT_OK)
33         goto error;
34
35     neat_start_event_loop(ctx, NEAT_RUN_DEFAULT);
36
37 error:
38     if (ctx)
39         neat_free_ctx(ctx);
40
41     return rc;
42 }
```

Listing 1: Setting up a listening socket with NEAT.

```
1 int
2 setup_listen_socket(const char* port, int file_count, char* file_names[])
3 {
4     struct addrinfo hints;
5     struct addrinfo *result, *rp;
6     int sfd, s;
7
8     memset(&hints, 0, sizeof(struct addrinfo));
9     hints.ai_family = AF_UNSPEC;
10    hints.ai_socktype = SOCK_STREAM;
11    hints.ai_flags = AI_PASSIVE;
12    hints.ai_protocol = IPPROTO_SCTP;
13    hints.ai_canonname = NULL;
14    hints.ai_addr = NULL;
15    hints.ai_next = NULL;
16
17    s = getaddrinfo(NULL, port, &hints, &result);
18    if (s != 0) {
19        fprintf(stderr, "getaddrinfo: %s\n", gai_strerror(s));
20        return -1;
21    }
22
23    for (rp = result; rp != NULL; rp = rp->ai_next) {
24        sfd = socket(rp->ai_family, rp->ai_socktype,
25                    rp->ai_protocol);
26        if (sfd == -1)
27            continue;
28
29        if (bind(sfd, rp->ai_addr, rp->ai_addrlen) == 0)
30            break; /* Success */
31
32        close(sfd);
```

```
33     }
34
35     if (rp == NULL) {
36         fprintf(stderr, "Could not bind\n");
37         return -1;
38     }
39
40     freeaddrinfo(result);
41
42     struct sctp_initmsg initmsg;
43     memset(&initmsg, 0, sizeof(initmsg));
44
45     initmsg.sinit_num_ostreams = file_count;
46
47     if (setsockopt(sfd, IPPROTO_SCTP, SCTP_INITMSG, &initmsg, sizeof(initmsg)) < 0)
48         return -1;
49
50     if (listen(sfd, SOMAXCONN) < 0) {
51         return -1;
52     }
53
54     listen_fd = sfd;
55     files = file_names;
56     no_of_files = file_count;
57
58     return 0;
59 }
```

Listing 2: Setting up a listening socket with the socket API.

## 2.3 Python bindings for the NEAT User API

Language bindings for the Python programming language have been generated, as a part of the continued development of the NEAT User API. This has not only been a feature request by potential external users, but it also serves to demonstrate flexibility with regard to programming paradigms. This section discusses the motivation and stumbling blocks, as well as the results achieved.

### 2.3.1 Motivation for NEAT Python bindings

The NEAT library is envisaged to become a core part of the future Internet architecture. In order to reach that goal, it has to support multiple platforms and programming environments. Thus, on the one hand, the NEAT library should run on several different Operating Systems. To this end, NEAT already supports FreeBSD, Linux, OS X and NetBSD.

On the other hand, however, it is also important to support more than just a single programming language. For comparison, the legacy Berkeley socket API is by now accessible in virtually every programming language which is productively used. This enables a multitude of interoperating pieces of software, tied together by the common grounds of a TCP/UDP transport system. It allows a client written in one language to connect and communicate with a server written in a totally different language.

In the long term, a similar flexibility should at least be possible for the NEAT System. Although NEAT itself is for a large part written in the C programming language, the option to use whatever language application programmers are comfortable with can substantially increase adaption rate and long-term viability and relevance of the NEAT library. In order to demonstrate that such a flexibility is indeed possible, bindings for an additional language have been generated for the NEAT library. The Python programming language has been chosen to serve as an example, demonstrating the basic flexibility of adapting NEAT concepts to a new programming environment.

Python is a scripting language that is both easy to learn and well known to many programmers. It is used in many high-profile software projects, such as Dropbox, YouTube, and Wikipedia. Its low barrier of entry makes it a favourable target for programming tutorials. For these reasons, it has been chosen as a target language for NEAT.

### 2.3.2 Using SWIG as solution approach

Making a C library usable from another language can be a tedious process of writing glue code. Luckily, there exist a number of frameworks that can help with this task. For interconnection with the Python programming language, the following options are available and have been considered:

- **Boost.python** (<http://www.boost.org/libs/python/>): This library is part of the larger Boost package, a set of C++ extension libraries. It serves well to connect C++ and Python and has also some advantages for connecting C with Python.
- **SIP** (<https://www.riverbankcomputing.com/software/sip/>): The SIP library (not to be confused with the Session Initiation Protocol, SIP) is another library that has been developed to simplify the creation of language bindings for Python. It has been used to create the PyQt library – Python language bindings for the Qt GUI toolkit. Like Boost.python, it is focused on Python only.
- **SWIG** (<http://www.swig.org/>): The “Simplified Wrapper and Interface Generator” is an interface compiler that can generate language bindings by parsing C or C++ header files and creating the necessary wrapper code automatically. It is not limited to just Python, but instead supports a large number of popular languages such as Perl, Java, or C#.

Due to its advantages, the SWIG framework was selected for the implementation of NEAT language bindings. It is a proven framework that has been used in many open-source projects, such as Subversion<sup>5</sup>, PyOpenGL<sup>6</sup>, or Xpian<sup>7</sup>. It supports creation of language bindings for a large number of well-known programming languages. Due to time and resource constraints, only bindings for Python are implemented in this project, but the use of SWIG will make it easy in any potential future development to integrate other languages as well.

The overall goal of this development effort was not to create a full-blown, production-ready interface. Rather, we wanted to demonstrate that it is possible with reasonable effort to make the NEAT library usable outside of a C/C++ scope without exposing the programmer of the target language to too many implementation details. More specifically, the following two goals were aimed for:

- To remain as close to the original interface as possible. Ideally, the programmer should be able to figure out the Python API without needing additional documentation. Additional helper functions that are not documented in the original C-based API will be avoided.

---

<sup>5</sup><https://subversion.apache.org/>

<sup>6</sup><http://pyopengl.sourceforge.net/>

<sup>7</sup><https://xpian.org/>

- To allow the programmer to write code in a “pythonic” way. Python programmers are used to a particular programming style, and the API should strive not to present any unexpected surprises.

A basic test case for this can be constructed by getting a Python server to talk to a C client, thereby demonstrating interoperability between two different programming languages. This approach was chosen here. To this end, the “minimal server” example shown in the NEAT online tutorial (§ 2.1.5) has been translated to Python.

### 2.3.3 Challenges and result: NEAT communication between C and Python

The generation of Python bindings with the SWIG framework is for a large part easy to accomplish. A SWIG definition file `neat_swig.i` has been written, based upon which SWIG can generate most of the necessary glue code itself. The translation of the example code into Python is then straightforward. To demonstrate, here the `on_writable` function used in the example code of `minimal_server` is used to compare the differences between its C and its Python implementation:

```
1 static neat_error_code
2 on_writable(struct neat_flow_operations *opCB)
3 {
4     const unsigned char message[] = "Hello, this is NEAT!";
5     neat_write(opCB->ctx, opCB->flow, message, 20, NULL, 0);
6     return NEAT_OK;
7 }
```

Listing 3: The `on_writable` function in C.

```
1 def on_writable(ops):
2     message = "Hello, this is NEAT!"
3     neat_write(ops.ctx, ops.flow, message, 20, None, 0)
4     return NEAT_OK
```

Listing 4: The `on_writable` function in Python.

The particular issue of callback functions makes the Python bindings more complicated. The NEAT library uses callback functions extensively to provide an event-oriented programming approach for the API users. These callbacks are defined in the `neat_flow_operations` struct in `neat.h` (cf. Appendix B). This problem is highlighted clearly in the `on_connected` function of the `minimal_server` example, which is shown here again in its C and Python implementations:

```
1 static neat_error_code
2 on_connected(struct neat_flow_operations *opCB)
3 {
4     opCB->on_writable = on_writable;
5     opCB->on_all_written = on_all_written;
6     neat_set_operations(opCB->ctx, opCB->flow, opCB);
7
8     return NEAT_OK;
9 }
```

Listing 5: The `on_connected` function in C.

```
1 def on_connected(ops):
2     ops.on_writable = on_writable
3     ops.on_all_written = on_all_written
4     neat_set_operations(ops.ctx, ops.flow, ops)
5     return NEAT_OK
```

Listing 6: The `on_connected` function in Python.

While, at first glance, the translation to Python appears to be straightforward, the concept of a Python program initializing a C library and asking it to call a Python function, again, make the implementation complicated. In line 2 of the Python function this means that the `ops.on_writable` on the left side of the assignment is a C construct, whereas the `on_writable` on the right side of the assignment is a Python function object. Clearly, intelligent type conversion, which remains transparent to the programmer, is necessary in this assignment to support such an intuitive programming style.

Although the SWIG library cannot solve this problem on its own, it still provides appropriate mechanisms (“typemaps”) to handle such situations. This has been used here to avoid having to change the API (e.g., by introducing a new converter method which the programmer would have to call manually). Instead, a type converter was defined for the SWIG library, which keeps an internal list of Python callback objects defined by the user, and uses C dispatcher functions to convert arguments and call the right Python function in time. To demonstrate, part of the code is shown here (the full code is available in `neat_swig.i`):

```
1 static struct {
2     PyObject *on_connected;
3     PyObject *on_error;
4     PyObject *on_readable;
5     PyObject *on_writable;
6     PyObject *on_all_written;
7     PyObject *on_network_status_changed;
8     PyObject *on_aborted;
9     PyObject *on_timeout;
10    PyObject *on_close;
11
12    PyObject *on_send_failure;
13    PyObject *on_slowdown;
14    PyObject *on_rate_hint;
15 } py_callbacks ;
16
17 static neat_error_code dispatch_fx(struct neat_flow_operations *ops, PyObject *
    pyfunc) {
18     PyObject *pyops = SWIG_NewPointerObj(SWIG_as_voidptr(ops),
        SWIGTYPE_p_neat_flow_operations, 0 | 0 );
19     PyObject *res = PyObject_CallFunctionObjArgs(pyfunc, pyops, NULL);
20     unsigned long val = PyLong_AsUnsignedLong(res);
21     return (neat_error_code)(val);
22 }
23
24 static neat_error_code disp_on_connected(struct neat_flow_operations *ops) {
25     return dispatch_fx(ops, py_callbacks.on_connected);
```



26 }

Listing 7: Wrapper code for callbacks in `neat_swig.i`.

Here, the `py_callbacks` structure provides an internal look-up table of user-defined Python callbacks. The `dispatch_fx` function translates the `neat_flow_operations` argument into a Python object, hands it to the specified Python function object, which is then called. Error codes are then translated back to C and returned to the caller. Finally, the `disp_on_connected` function serves as a simple wrapper, feeding the right Python function object to the actual dispatch function. This mechanism now allows the user to define appropriate NEAT callback functions in Python and plug them into the NEAT framework.

### 2.3.4 Summary

With the above issues solved, the Python bindings are available for programmers to use. For demonstration purposes, the `minimal_server` example has been completely translated into Python and is available in the NEAT Github source code repository. While the feature is still considered experimental and, thus, has to be activated deliberately, still the barrier of entry for NEAT is lowered further. With a first implementation operational, programmers now have the chance to test and improve the availability of NEAT features in Python.

## 3 Core Transport Functions

**Note:** The descriptions of NEAT components presented in this section **replace** those in Deliverable D2.2, reflecting the status of the implementation at the end of Work Package 2 activities.

This section presents a detailed description of the transport functionalities required to realise the NEAT core transport system based on the components introduced in Section 1.3. Each component is described in detail and Transport Service Features they provide are specified when applicable. Indicative examples of their operation are given where relevant, and relationships among these building blocks are identified. This document does not aim to provide all the implementation details, rather, it intends to present the design choices that have been made. For a more comprehensive reference of the NEAT User API, please refer to Appendix B (also available online at: <http://neat.readthedocs.io/en/latest/index.html>). Where appropriate, snippets of sample code in C are used for better presentation.

### 3.1 NEAT Framework Components

To run a NEAT System a minimum set of basic building blocks has to be implemented, comprising the NEAT Framework components.

This translates into being able to create a NEAT Flow and connect to a host using a domain name address, as well as the ability to translate the functionalities behind the NEAT User API into appropriate function calls, e.g., to different protocols, mechanisms, etc.

Setting up a NEAT Flow can be done by using an event-based, user-space *NEAT library* that implements a callback-based API (NEAT API Framework). Once a NEAT Flow is initialised, it will contain a

structure that keeps all of its relevant information during its lifetime (NEAT Flow Endpoint). A NEAT Flow can be assigned to one or more domain names as well as IP addresses (Connect to a name). The functionalities behind the NEAT User API requested for the initialised NEAT Flow require code that “glues” together different components (NEAT Logic). This is not a monolithic chunk of code separated from other components, but rather code that is scattered throughout other components as well as the NEAT User API. Finally, gathering statistics and information about the operation of the system is necessary for diagnostics and performance monitoring (NEAT Flow Endpoint Statistics). The operation of each of these components is presented in the rest of this subsection.

### 3.1.1 NEAT Flow Endpoint

The NEAT Flow Endpoint is a NEAT structure that has a role similar to that of the Transmission Control Block (TCB) in the context of TCP [15]. This provides access to the additional flow information that extends the transport components and realises a transport-independent interface. Through the NEAT Logic parts of this information are used by most of the other building blocks (e.g., Policy Manager, Happy Eyeballs, Security, etc.).

A NEAT Flow Endpoint `neat_flow` structure corresponds to a single operating system socket or to a stream if multistreaming is used (only for multistreaming capable protocols, like SCTP) and keeps the information about the socket or stream that is relevant during the NEAT Flow’s lifetime (Listing 8). This includes:

- A structure that holds information about the underlying OS socket as well as the pointer to the OS file descriptor or a user-space SCTP socket. This information includes: socket family (e.g. `AF_INET`, `AF_INET6`), socket type (e.g., `SOCK_DGRAM`, `SOCK_STREAM`), protocols (e.g., `IPPROTO_UDP`, `IPPROTO_TCP`, `IPPROTO_SCTP`)
- Remote peer domain name.
- Server certificates and key in PEM format [12, 27, 28, 31].
- Remote port.
- Remote socket address.
- Local addresses.

```
1 struct neat_flow
2 {
3     struct neat_pollable_socket *socket;
4     TAILQ_HEAD(neat_listen_socket_head, neat_pollable_socket) listen_sockets;
5
6     const char *name;
7     char *server_pem;
8     char *key_pem;
9     uint16_t port;
10    const struct sockaddr *sockAddr; // raw unowned pointer into resolver_results
11    json_t *user_ips;
```

Listing 8: NEAT Flow Endpoint structure.

It also contains information relevant to the NEAT System operations (e.g., DNS resolver outcome), flow statistics and desired properties set by an application (e.g., QoS, ECN, etc.) and includes (Listing 9):

- Address resolver results, i.e., DNS results.
- Flow statistics: this includes the path Maximum Transmission Unit (MTU), slow-start threshold, round-trip time (RTT), etc.
- Requested properties for the NEAT Flow: these are properties that an application has specified when the flow has been created via the NEAT User API. They are in JSON format (see Listing 24). The Policy Manager translates these properties into a list of candidates.
- Candidate list: it is used by the Happy Eyeballs component (§ 3.3.1) to create concrete socket configuration probes.
- Connection attempt count: the number of different configurations the Happy Eyeballs component has tried.
- QoS: the *abstract* Quality of Service (QoS) parameter to be used with this flow (see Table 2).
- ECN: a value to indicate if use of Explicit Congestion Notification (ECN) should be attempted.
- Group: the group ID that this flow belongs to. This is used for coupled congestion control (§ 3.2.3).
- Flow priority: the priority of the given NEAT flow relative to the other flows in the same group. Must be between 0.1 and 1.0.
- Congestion control: The congestion control algorithm to be used for this flow.

```
12 struct neat_resolver_results *resolver_results;
13 struct neat_flow_statistics flow_stats;
14 json_t *properties;
15 struct neat_he_candidates *candidate_list;
16 uint8_t heConnectAttemptCount;
17
18 uint8_t qos;
19 uint8_t ecn;
20
21 uint32_t group;
22 float priority;
23
24 const char *cc_algorithm;
```

Listing 9: NEAT Flow Endpoint structure (continued).

The NEAT System provides a callback-based API to the application. Pointers to the callback functions are kept in the NEAT Flow Endpoint structure (Listing 10). Their usage will be clarified in § 3.1.2.

```
25 struct neat_flow_operations *operations;
```

Listing 10: NEAT Flow Endpoint structure (continued).

The flow endpoint structure is used as the main hub for communication with the underlying socket. Therefore it contains the pointers to the NEAT base loop structure `neat_ctx`, functions for accessing the underlying socket, and NEAT Flow internal flags and buffers (the write buffer facilitates preservation of message boundaries if the selected transport protocol is message-based), see Listing 11.

```
26 // NEAT base loop structure:
27 struct neat_ctx *ctx;
28
29 // NEAT flow state.
30 uint8_t state;
31
32 // Functions for accessing the underlying socket:
33 neat_read_impl readfx;
34 neat_write_impl writefx;
35 neat_accept_impl acceptfx;
36 neat_connect_impl connectfx;
37 neat_close_impl closefx;
38 neat_listen_impl listenfx;
39 neat_shutdown_impl shutdownfx;
40
41 // NEAT internal flags:
42 unsigned int hefirstConnect      : 1;
43 unsigned int firstWritePending  : 1;
44 unsigned int acceptPending      : 1;
45 unsigned int isPolling          : 1;
46 unsigned int everConnected      : 1;
47 unsigned int isDraining         : 1;
48 unsigned int isServer           : 1;
49 unsigned int isSCTPMultihoming  : 1;
50 unsigned int security_needed    : 1;
51 unsigned int isSCTPIData        : 1;
52 unsigned int isClosing          : 1;
53 unsigned int notifyDrainPending : 1;
54 unsigned int preserveMessageBoundaries : 1;
55 unsigned int eofSeen            : 1;
56 unsigned int skipCertVerification : 1;
57
58 // Write buffer:
59 struct neat_message_queue_head bufferedMessages;
60
61 uv_poll_cb callback_fx;
62 struct neat_iofilter *iofilters;
```

Listing 11: NEAT Flow Endpoint structure (continued).

SCTP needs additional internal flow states and variables shown in Listing 12. This includes additional information needed if user-space SCTP or multistreaming are configured.

```
63 // The memory buffer for reading. Used of SCTP reassembly.
64 unsigned char *readBuffer; // memory for read buffer
65 size_t readBufferSize; // amount of received data
```

```
66     size_t          readBufferAllocation; // size of buffered allocation
67     int readBufferMsgComplete; // it contains a complete user message
68
69     unsigned int streams_requested;
70
71 #if defined(USR_SCTP_SUPPORT)
72     neat_accept_usr_sctp_impl accept_usr_sctp_fx;
73 #endif
74
75
76 #ifndef SCTP_MULTISTREAMING
77     unsigned int multistream_check : 1;
78     unsigned int multistream_shutdown : 1;
79     unsigned int multistream_reset_in : 1;
80     unsigned int multistream_reset_out : 1;
81
82     uv_timer_t *multistream_timer;
83     uint16_t multistream_id;
84     LIST_ENTRY(neat_flow) multistream_next_flow;
85
86     struct neat_read_queue_head multistream_read_queue;
87     size_t multistream_read_queue_size;
88
89     neat_flow_states multistream_state;
90 #endif
```

Listing 12: NEAT Flow Endpoint structure (continued).

The NEAT System supports TCP Fast Open [16] and Listing 13 shows states needed by this feature, e.g., a buffer that holds data that should be sent on a SYN packet or after the connection is established for protocols that do not support this feature, etc.

```
92     unsigned char *tfoBuffer; // buffer for tcp fast open.
93     size_t tfoBufferWritten; // amount of data in the tfo buffer.
94     struct neat_tlv *tfoOptions;
95     uint32_t tfoOptionsCount;
96 }
```

Listing 13: NEAT Flow Endpoint structure (continued).

**Some examples of the operation:** Listing 14 shows an example of how the NEAT Flow Endpoint structure is used. After DNS is resolved a list of candidates is handed over to the Happy Eyeballs component and this component will try to establish connections using some or all of the given candidates. After the first candidate is connected, the NEAT System has all socket parameters and these parameters are stored in the `neat_flow` structure. The other candidates are released. This is done in the `he_connected_cb` callback function. For simplicity, some parts of the code are omitted.

The `he_connected_cb` function is called when a connection is established or an error occurs. The function parameter `uv_poll_t` holds a pointer to the candidate, i.e., a `neat_he_candidate` structure. The `neat_he_candidate` structure holds a pointer to the corresponding `neat_flow` structure.

When the first candidate is successfully connected, the `neat_flow` structure is filled with the necessary information (lines 34-47). When the other candidates are connected, the corresponding socket is closed (lines 50-63).

```
1 static void
2 he_connected_cb(uv_poll_t *handle, int status, int events)
3 {
4     struct neat_he_candidate *candidate = handle->data;
5     struct neat_flow *flow = candidate->pollable_socket->flow;
6     struct neat_he_candidates *candidate_list = flow->candidate_list;
7
8     assert(candidate);
9     assert(candidate->pollable_socket);
10    assert(flow);
11
12    int so_error = 0;
13    unsigned int len = sizeof(so_error);
14    if (getsockopt(candidate->pollable_socket->fd, SOL_SOCKET, SO_ERROR, &so_error,
15                &len) < 0) {
16        neat_log(NEAT_LOG_DEBUG, "Call to getsockopt for fd %d failed: %s",
17                candidate->pollable_socket->fd, strerror(errno));
18
19        uv_poll_stop(handle);
20        uv_close((uv_handle_t*)handle, free_he_handle_cb);
21
22        neat_io_error(candidate->ctx, flow, NEAT_ERROR_INTERNAL);
23        return;
24    }
25    status = so_error;
26    neat_log(NEAT_LOG_DEBUG, "%s - Connection status: %d", __func__, status);
27
28    if (flow->hefirstConnect && (status == 0)) {
29        // This is the fastest candidate.
30        // Change the flow internal state after one of the candidates was
31        // successfully connected.
32        flow->hefirstConnect = 0;
33
34        assert(flow->socket);
35
36        flow->socket->fd = candidate->pollable_socket->fd;
37        flow->socket->flow = flow;
38        flow->socket->handle = handle;
39        flow->socket->handle->data = flow->socket;
40        flow->socket->family = candidate->pollable_socket->family;
41        flow->socket->type = candidate->pollable_socket->type;
42        flow->socket->stack = candidate->pollable_socket->stack;
43        flow->socket->write_size = candidate->pollable_socket->write_size;
44    }
45    ;
46}
```

```
42     flow->socket->write_limit           = candidate->pollable_socket->
        write_limit;
43     flow->socket->read_size             = candidate->pollable_socket->read_size;
44     flow->socket->sctp_explicit_eor     = candidate->pollable_socket->
        sctp_explicit_eor;
45
46     flow->everConnected = 1;
47     flow->isPolling = 1;
48
49 } else {
50     // This is not the first candidate so close the socket.
51     neat_log(NEAT_LOG_DEBUG, "%s - NOT first connect", __func__);
52
53     uv_poll_stop(handle);
54     uv_close((uv_handle_t*)handle, free_he_handle_cb);
55
56     neat_log(NEAT_LOG_DEBUG, "%s:Release candidate", __func__);
57     TAILQ_REMOVE(candidate_list, candidate, next);
58     free(candidate->pollable_socket->dst_address);
59     free(candidate->pollable_socket->src_address);
60     free(candidate->pollable_socket);
61     free(candidate->if_name);
62     json_decref(candidate->properties);
63     free(candidate);
64 }
65 }
```

Listing 14: Use of the NEAT Flow Endpoint structure.

**Provided Transport Service Feature(s):** This building block is part of the most basic functionality of the NEAT System and does not relate to any specific application requirement.

**Related building blocks:**

- NEAT Logic (§ 3.1.3).

### 3.1.2 NEAT API Framework (callback)

The NEAT System implements a callback-based API based on the libuv [3] library to provide portable asynchronous I/O across multiple platforms. The base of the NEAT System is an event loop that needs to be initialised before any NEAT functionality can be accessed. NEAT uses libuv [3] as an event library. Once the NEAT base structure has started, an application can request a connection (create NEAT Flow), communicate over it (write data to the NEAT Flow and read received data from the NEAT Flow) and register callback functions that will be executed upon the occurrence of certain events.

The NEAT System offers several basic functions for accessing a NEAT Flow (Listing 15):

- Creating and closing a NEAT Flow, i.e., a NEAT Flow Endpoint structure (this structure is described in § 3.1.1).

- Setting callback functions.
- Opening a connection to a remote host: this function starts the NEAT Logic for selecting and creating the most adequate Transport Service. The outcome of this asynchronous call can be an `on_connected` event in case of success, or an `on_error` event if a failure occurs. Besides a remote host name and a port number, this function can take some additional parameters such as the flow priority, a local address to be used, requesting multihoming, etc. These additional parameters are optional. The NEAT System supports the TCP Fast Open extension [16] that sends data on the TCP SYN packet. An extended version of the `neat_open` function, called `neat_open_with_data`, accepts the data to be sent on a SYN packet as a function parameter as well (`neat_open_with_data` corresponds to the `OPEN_WITH_EARLY_DATA` primitive from D1.3 [41]).
- Read from, and write to, a NEAT Flow and through it from/to the underlying socket: the return values correspond to the return values of the operating system function calls, or of the TLS/DTLS function calls if secure communication is established.
- Functions for getting and setting properties: using these functions, an application can set or change transport requirements (e.g., reliable transport, low latency, etc.). The properties are in a JSON format, presented in Listing 24.
- Opening a NEAT Flow for listening (i.e., server-side socket): this triggers an `on_connected` event if a new connection request is received or an `on_error` event in the case of an error. Besides a port number, a local address to be used and the maximal number of streams can be specified. These parameters are optional.
- Getting flow statistics (e.g., round-trip time, slow-start threshold, etc.).
- Querying the flow for local and remote peer addresses.
- Setting a user specified timeout, i.e., `TCP_USER_TIMEOUT`.
- Setting some parameters that are specific to certain protocols, e.g., `neat_set_primary_dest` only used with SCTP and `neat_set_checksum_coverage` only used for UDP and UDP-Lite.
- Setting server certificates (i.e., `neat_secure_identity`).
- Setting additional desired parameters, such as QoS, ECN and the `TCP_NOTSENT_LOWAT` parameter to limit amount of unsent data in a TCP socket.

```
1 struct neat_flow *neat_new_flow(struct neat_ctx *ctx);
2
3 neat_error_code neat_shutdown(struct neat_ctx *ctx, struct neat_flow *flow);
4 neat_error_code neat_close(struct neat_ctx *ctx, struct neat_flow *flow);
5 neat_error_code neat_abort(struct neat_ctx *ctx, struct neat_flow *flow);
6
7 neat_error_code neat_set_operations(struct neat_ctx *ctx, struct neat_flow *flow,
8                                   struct neat_flow_operations *ops);
9 neat_error_code neat_open(struct neat_ctx *ctx, struct neat_flow *flow,
10                           const char *name, const char *port,
11                           struct neat_tlv optional[], unsigned int opt_count);
```



```
12 neat_error_code neat_open_with_data(struct neat_ctx *mgr, struct neat_flow *flow,
13                                     const char *name, uint16_t port,
14                                     struct neat_tlv optional[], unsigned int opt_count,
15                                     const unsigned char *buffer, uint32_t amt, uint32_t *
                                         written);
16 neat_error_code neat_read(struct neat_ctx *ctx, struct neat_flow *flow,
17                             unsigned char *buffer, uint32_t amt, uint32_t *actualAmt,
18                             struct neat_tlv optional[], unsigned int opt_count);
19 neat_error_code neat_write(struct neat_ctx *ctx, struct neat_flow *flow,
20                             const unsigned char *buffer, uint32_t amt,
21                             struct neat_tlv optional[], unsigned int opt_count);
22
23 neat_error_code neat_get_property(struct neat_ctx *ctx, struct neat_flow *flow,
24                                   const char* name, void *ptr, size_t *size);
25 neat_error_code neat_set_property(struct neat_ctx *ctx, struct neat_flow *flow,
26                                   const char* properties);
27
28 neat_error_code neat_accept(struct neat_ctx *ctx, struct neat_flow *flow,
29                             uint16_t port, struct neat_tlv optional[],
30                             unsigned int opt_count);
31
32 neat_error_code neat_get_stats(struct neat_ctx *ctx, char **neat_stats);
33
34 int neat_getlpadrs(struct neat_ctx *ctx, struct neat_flow *flow,
35                   struct sockaddr** adrs, const int local);
36 void neat_freelpadrs(struct sockaddr* adrs);
37
38 neat_error_code neat_change_timeout(struct neat_ctx *ctx, struct neat_flow *flow,
39                                     unsigned int seconds);
40
41 neat_error_code neat_set_primary_dest(struct neat_ctx *ctx, struct neat_flow *flow,
42                                       const char *name);
43 neat_error_code neat_set_checksum_coverage(struct neat_ctx *ctx,
44                                             struct neat_flow *flow,
45                                             unsigned int send_coverage,
46                                             unsigned int receive_coverage);
47
48 neat_error_code neat_secure_identity(struct neat_ctx *ctx, struct neat_flow *flow,
49                                       const char *filename, int pemType);
50
51 neat_error_code neat_set_qos(struct neat_ctx *ctx, struct neat_flow *flow,
52                               uint8_t qos);
53 neat_error_code neat_set_ecn(struct neat_ctx *ctx, struct neat_flow *flow,
54                               uint8_t ecn);
55
56 neat_error_code neat_set_low_watermark(struct neat_ctx *ctx, struct neat_flow *flow,
57                                         uint32_t watermark);
```

Listing 15: NEAT API functions.

The NEAT API offers multiple run-time events that call the corresponding functions if registered. The set of events that are triggered are given in Listing 16.

```
1 neat_flow_operations_fx on_connected;
2 neat_flow_operations_fx on_error;
3 neat_flow_operations_fx on_readable;
4 neat_flow_operations_fx on_writable;
5 neat_flow_operations_fx on_all_written;
6 neat_flow_operations_fx on_network_status_changed;
7 neat_flow_operations_fx on_aborted;
8 neat_flow_operations_fx on_timeout;
9 neat_flow_operations_fx on_close;
10 neat_cb_send_failure_t on_send_failure;
11 neat_cb_flow_slowdown_t on_slowdown;
12 neat_cb_flow_rate_hint_t on_rate_hint;
```

Listing 16: NEAT callback functions.

In the following, we present a simple illustration of the NEAT Flow connection establishment and maintenance. Each application needs to initialise the base NEAT structure and before a communication can start its event loop needs to be started. An application creates a NEAT Flow Endpoint for each connection. The application specifies its requirements by utilising the `neat_set_property` and `neat_get_property` functions of the NEAT User API called for each individual NEAT Flow Endpoint. When the application requirements are specified, the connection establishment can be requested by calling the `neat_open` function or the `neat_open_with_data` function (corresponding to the OPEN primitive and the OPEN\_WITH\_EARLY\_DATA primitive from D1.3 [41]). The function takes a host name and a port number as parameters and maybe some optional parameters. On the server side, the `neat_accept` function (corresponding to the ACCEPT primitive from D1.3 [41]) will be called with parameters: port number and local address the socket should be listening to, and optionally additional parameters.

These two function calls will trigger a set of actions inside the NEAT System. The specified application requirements will be used by the NEAT Logic and the corresponding building blocks (e.g., Policy Manager and Happy Eyeballs) to select and probe selected socket configurations. If a socket that satisfies the application requirements has been successfully connected, an `on_connected` callback function, if registered, will be invoked. Otherwise the `on_error` callback function will be invoked.

The application can register callback functions for `on_socket_readable` and `on_socket_writable` events which will translate into poll parameters for the underlying socket. The functions will be executed if the corresponding event applies.

The NEAT System buffers data that needs to be written. This is necessary to facilitate preservation of message boundaries if the selected transport protocol is message-based. The `on_all_written` event is triggered when all buffered data is written out to the OS socket.

In case of an error (e.g., NEAT internal error, socket error, socket being closed, etc.) an `on_error` event callback function will be invoked.

**Some examples of the operation:** Listing 17 presents a simple example of setting callback functions and waiting for a callback function to be called. In function `main` a NEAT base loop structure and a NEAT Flow are created (lines 21 and 28). Callback functions `on_error` and `on_connected` are set in lines 59–60. `neat_open` is called in line 70 and if it does not return an error the NEAT loop will be

started (line 71). In case of an error the `on_error` function will be invoked, otherwise `on_connected` will be invoked which sets callback functions for `on_all_written` and `on_readable`.

```
1
2 static struct neat_flow_operations ops;
3
4 /*
5     Error handler
6 */
7 static neat_error_code on_error(struct neat_flow_operations *opCB)
8 {
9     exit(EXIT_FAILURE);
10 }
11
12 static neat_error_code on_connected(struct neat_flow_operations *opCB)
13 {
14     opCB->on_all_written = on_all_written;
15     opCB->on_readable = on_readable;
16     return NEAT_OK;
17 }
18
19 int main(int argc, char *argv[])
20 {
21     if ((ctx = neat_init_ctx()) == NULL) {
22         debug_error("could not initialise context");
23         result = EXIT_FAILURE;
24         goto cleanup;
25     }
26
27     // new neat flow
28     if ((flow = neat_new_flow(ctx)) == NULL) {
29         debug_error("neat_new_flow");
30         result = EXIT_FAILURE;
31         goto cleanup;
32     }
33
34     char *config_property = "{\
35         \"transport\": [\
36             {\
37                 \"value\": \"SCTP\", \
38                 \"precedence\": 1 \
39             }, \
40             {\
41                 \"value\": \"SCTP/UDP\", \
42                 \"precedence\": 1 \
43             }, \
44             {\
45                 \"value\": \"TCP\", \
46                 \"precedence\": 1 \
```

```
47         }\
48     ]\
49 }";
50
51 // set properties
52 if (neat_set_property(ctx, flow, config_property)) {
53     fprintf(stderr, "%s - error: neat_set_property\n", __func__);
54     result = EXIT_FAILURE;
55     goto cleanup;
56 }
57
58 // set callbacks
59 ops.on_connected = on_connected;
60 ops.on_error = on_error;
61
62 if (neat_set_operations(ctx, flow, &ops)) {
63     debug_error("neat_set_operations");
64     result = EXIT_FAILURE;
65     goto cleanup;
66 }
67
68 // wait for on_connected or on_error to be invoked
69 // The last 2 arguments are the remote server name and port number.
70 if (neat_open(ctx, flow, argv[argc - 2], argv[argc - 1]) == NEAT_OK) {
71     neat_start_event_loop(ctx, NEAT_RUN_DEFAULT);
72 } else {
73     debug_error("neat_open");
74     result = EXIT_FAILURE;
75     goto cleanup;
76 }
77
78 ...
79 }
```

Listing 17: NEAT API Framework example.

**Provided Transport Service Feature(s):** This building block is part of the most basic functionality of the NEAT System and it does not relate to any specific Transport Service Feature.

**Related building blocks:**

- NEAT Logic (§ 3.1.3).

### 3.1.3 NEAT Logic

The NEAT Logic component is at the core of the NEAT System as part of the NEAT Framework components and is responsible for providing functionalities behind the NEAT User API. It orchestrates and “glues” together different components. The NEAT Logic is not a monolithic piece of code separated

from other components, but its code is scattered throughout other components as well as the NEAT User API.

Requests made via the NEAT User API are translated into function calls to the Policy Manager or other building blocks; for instance, calls to select the transport protocols to be instantiated, or calls to Handover and Signalling after receiving a set of candidates from the Policy Manager. Transport protocols are configured via the relevant NEAT Transport components. The NEAT Logic dispatches different decisions returned by the Policy Manager by translating them into certain function calls related to NEAT components—e.g., by calling the Happy Eyeballs function(s) for SCTP/TCP or IPv6/IPv4. In simpler terms, it glues different building blocks of the NEAT System together and makes them operational in one uniform system.

The NEAT Logic also maps the primitives and events exposed to the application (as described in Deliverable D1.3 [41]) to the primitives provided by each transport protocol.

A detailed sequence of the NEAT Logic operation with regards to the connection setup in NEAT is presented in Section 1.4 and therefore we avoid repeating it in here for the sake of brevity.

**Some examples of the operation:** An example of the NEAT Logic operation is when the `neat_open` call, corresponding to the `OPEN` primitive (see § 2.1.2 of D1.3 [41]), is used to open a NEAT Flow. The `OPEN` primitive does not specify any specific transport protocol. After processing the optional arguments passed in the `neat_tlv` struct (e.g., priority, flow group or congestion control), and initialising the NEAT Flow's address name and port, it would send the properties to the Policy Manager using the `send_properties_to_pm()` function<sup>8</sup>. This corresponds to the first query to the Policy Manager as depicted in Figure 4 and explained in Section 1.4.

```
1 neat_error_code
2 neat_open(neat_ctx *ctx, neat_flow *flow, const char *name, uint16_t port,
3           struct neat_tlv optional[], unsigned int opt_count)
4 {
5 ...
6
7     HANDLE_OPTIONAL_ARGUMENTS_START()
8         OPTIONAL_INTEGER(NEAT_TAG_STREAM_COUNT, stream_count)
9         OPTIONAL_INTEGER(NEAT_TAG_FLOW_GROUP, group)
10        OPTIONAL_FLOAT(NEAT_TAG_PRIORITY, priority)
11        OPTIONAL_STRING(NEAT_TAG_CC_ALGORITHM, cc_algorithm)
12    HANDLE_OPTIONAL_ARGUMENTS_END();
13
14 ...
15
16    flow->name = strdup(name);
17    if (flow->name == NULL)
18        return NEAT_ERROR_OUT_OF_MEMORY;
19    flow->port = port;
20    flow->group = group;
21    flow->priority = priority;
22
23 ...
```

---

<sup>8</sup>The flow properties should be already set by `neat_set_property()` prior to the `neat_open` function call.

```
24
25     send_properties_to_pm(ctx, flow);
26     return NEAT_OK;
27 }
```

Listing 18: NEAT open function.

**Provided Transport Service Feature(s):** There are no specific Transport Service Features associated to this building block. However, the operation of NEAT Logic is essential for other building blocks to provide their Transport Service Features.

**Related building blocks:**

- NEAT Flow Endpoint Statistics (§ 3.1.5).
- Middlebox Traversal (§ 3.2.2).
- NEAT Flow Endpoint (§ 3.1.1).
- NEAT API Framework (callback) (§ 3.1.2).
- Connect to a name (§ 3.1.4).
- Happy Eyeballs (§ 3.3.1).
- Security (§ 3.2.4).
- NEAT Policy Manager (via Policy Interface) (§ 3.4.1).

### 3.1.4 Connect to a name

The higher-level API offered by NEAT provides transport-independent name resolution. This approach not only avoids dependencies on specific network technology, but is also essential to support multiple active network interfaces. The NEAT system address resolver is named *Connect to a name*. NEAT will by default resolve domains using the public DNS servers provided by Google and OpenDNS (both v4 and v6). If a system-wide resolver file is found, then the list is extended with the servers contained in this file.

Currently we support resolving names using plain DNS (i.e., no DNSSEC), but the component is designed in such a way that extending the functionality is easy. One can for example imagine that a new name resolution scheme, protocol or technique will be introduced later on. An IP literal can also be provided, in which case no translation will be performed.

The resolver supports multi-homed hosts. When the NEAT resolve function is called, the query will by default be sent over all available (interface, address) tuples. Only A records are requested over IPv4 addresses, while AAAA records are requested over the available IPv6 addresses. A challenge on multi-homed hosts, which cannot be solved within NEAT, is to make sure that we use the correct interface/address when communicating with a DNS server. This is an issue when for example a host is connected to overlapping networks, or when the DNS servers acquired through DHCP have public IP addresses. There is currently no standardised syntax for how to specify which local interface/address should be used to communicate with a DNS server. We designed, implemented and shared a solution

for the popular open-source resolver *dnsmasq*<sup>9</sup>, but the discussions with the *dnsmasq* community have not progressed yet.

Interface/address tuples used by NEAT are stored in a list which is dynamically updated based on events generated by the OS. All major OS support mechanisms for generating events when addresses/network interfaces are added/removed. The mechanisms differ based on OS, so small shims are needed to support different operating systems. However, the core code (and the content of the list) is platform-independent.

The interface/address list is stored in the `neat_ctx` object introduced in § 3.1.1 and is available for use by all other building blocks. For example the address monitoring functionality, combined with an internal notification subsystem, will make reacting to interfaces going up/down more efficient across all blocks.

In summary, the following features are provided by the *Connect to a name* component:

- *Asynchronous DNS lookup*: name resolving will not block the calling application.
- *Address monitoring*: (interface, address) tuples are stored in a dynamically updated list, which is available to all building blocks.
- *Multi-homing support*: the resolver will resolve names using all (interface, address) tuples on a host by default.
- *Private network marking*: names resolving to internal addresses will be marked and can be easily filtered.

**Some examples of the operation:** To use the resolver, a developer has to create the NEAT context first using the `neat_init_ctx` call. Then, the `neat_resolver_init` function has to be called to set up the resolver. This function is passed two function pointers, one that will be called when the resolver finishes (or times out), and one called when the resolver can be released.

After these two functions have been called, it is simply a matter of calling `neat_getaddrinfo`. This function works very similar to the POSIX-compliant `getaddrinfo`. In other words, it is possible to limit which address family and transport protocol for the resolver to query/return. One change from the normal `getaddrinfo` is that returning multiple transport protocols is supported.

When the resolving is done (or has failed, e.g., due to a timeout), the provided callback function is called. If successful, this function is passed a list of all (interface, source address, destination address) tuples. For example these can be used by a developer to connect to the desired host, or the list will serve as input to the Happy Eyeballs component.

**Provided Transport Service Feature(s):** Connect to a name.

#### **Related building blocks:**

- Connect to a name has no dependencies on other building blocks except NEAT Logic (§ 3.1.3), but several building blocks may depend on the offered functionality. One example is Happy Eyeballs (§ 3.3.1) that must be provided with a set of source/destination addresses to probe for IPv4/IPv6 and transport-protocol connectivity.

<sup>9</sup><http://www.thekelleys.org.uk/dnsmasq/doc.html>

### 3.1.5 NEAT Flow Endpoint Statistics

NEAT not only automates important network decisions for applications, it can also help understand how these decisions were taken. By making this flow information available to the user in a consistent form, it eases the burden of identifying the root causes of any network problems.

The NEAT Flow Endpoint Statistics component is responsible for maintaining information about the current state of the NEAT System, and for gathering usage statistics for both the overall NEAT System and the respective NEAT Flows. This information resembles the information provided by `netstat` in a traditional socket stack, but at the NEAT Flow level of detail. It can be used for application-level diagnostic purposes and for allowing applications to monitor the performance of application flows (e.g., measuring the throughput of NEAT Flows) to take decisions based on provided detailed information.

The information provided by the NEAT Flow Endpoint Statistics building block can be divided into three sets: 1) *current NEAT state*, 2) *NEAT Flow statistics*, and 3) *NEAT System statistics*. The *current NEAT state* set of information provides a detailed view of the current state of the NEAT System. It provides a list of all NEAT Flows that are currently open on the NEAT System along with details about their configuration. The *NEAT Flow statistics* set of information contains historical statistics since the NEAT instance was created. The *NEAT System statistics* set of information contains system-wide usage statistics of the NEAT System. Some examples of these three categories are listed below:

- *Current NEAT state*: flow ID, flow creation time, local name, local transport address(es), destination name, destination transport address(es), send queue size, protocol state, transport parameters (e.g., Nagle, DSCP, timeout, etc.), flow properties, interface(s) in use.
- *NEAT Flow statistics*: number of bytes sent/received, number of messages sent/received, number of messages dropped locally, number of handovers, connection duration, creation time.
- *NEAT System statistics*: total number of bytes sent/received, total number of messages sent/received, total number of messages dropped locally, total number of opened/accepted/closed/aborted connections,

The number and content of the elements returned after a call to the statistics interface will vary between NEAT instances, and even within the lifetime of one NEAT context. Different operating systems, for example, support different transport options, meaning that the NEAT state will look different between operating systems. Also, the protocol used for a particular flow will be different between connections, meaning that a different set of available options will have to be presented each time the statistics module is called. Each network protocol and implementation also has differences in how many statistics are made available, leading to differences in what statistics can be provided by the NEAT framework depending on the OS it is running on and on the protocol stack used on that OS (e.g., userland SCTP will expose a different set of statistics than kernel SCTP). We have, so far, not focused on extending the native statistics support for each protocol implementation. We are rather trying to expose the most useful statistics depending on the available statistics provided by the protocols chosen by the NEAT framework. The variety of the data to be given to the caller makes it a necessity to use a format that supports such dynamic content. In line with the standards chosen for the PM, the statistics module therefore employs a JSON format supporting different elements in a nested structure of objects. Listing 19 shows an example of how statistics are presented in this format. An important consideration is to make sure that the JSON statistics elements are extended/evolved in a backwards



compatible way, such that applications will always be able to collect the statistics they expect in the required format.

```
1 {
2   "flow1": {
3     "flow number": 1,
4     "sock_protocol": 6,
5     "remote_host": "localhost",
6     "readSize": 87380,
7     "bytes sent": 2712,
8     "socket type": 1,
9     "tcpstats": {
10      "reordering": 3,
11      "retransmits": 0,
12      "rtt": 19,
13      "pmtu": 65535,
14      "rttvar": 8,
15      "ssthresh": 2147483647,
16      "rcv_ssthresh": 43690,
17      "snd_cwnd": 10
18    },
19    "port": 8080,
20    "writeSize": 16384,
21    "bytes received": 1337
22  },
23  "flow2": {
24    "flow number": 1,
25    "sock_protocol": 6,
26    "remote_host": "localhost",
27    "readSize": 87380,
28    "bytes sent": 6,
29    "socket type": 1,
30    "tcpstats": {
31      "reordering": 3,
32      "retransmits": 0,
33      "rtt": 19,
34      "pmtu": 65535,
35      "rttvar": 8,
36      "ssthresh": 2147483647,
37      "rcv_ssthresh": 43690,
38      "snd_cwnd": 10
39    },
40    "port": 8080,
41    "writeSize": 16384,
42    "bytes received": 0
43  },
44  "number of flows": 2,
45  "global bytes sent": 2718,
46  "global bytes received": 1337
```

47 }

Listing 19: Example of the NEAT Statistics JSON format.

Applications that want to take advantage of this additional information provided by the NEAT System (i.e., Class-4 applications in Figure 2) can access it through the Diagnostics and Statistics Interface. The scope of the information maintained by the NEAT Flow Endpoint Statistics is local to the application that uses a particular NEAT System instance and is maintained within the application context as long as the application is running.

To provide system administrators with a tool to collect network statistics across NEAT instances, the statistics will have to be collected as a CIB source (§ 3.4.3). A *neatstat* tool, further described in Deliverable D4.2 [14], can then connect to a dedicated socket of the PM to be given updated statistics for all NEAT contexts running on a machine.

**Some examples of the operation:** Listing 20 shows an example of how the statistics can be called by an application. In the example, the JSON string is simply printed to the standard output, but normally the application would parse the string, pick the relevant elements and process the content. Since the statistics are passed by reference, the reserved memory needs to be freed by the calling application. This is done for efficiency reasons, to save system resources.

```
1 static void
2 print_neat_stats(neat_ctx *mgr)
3 {
4     neat_error_code error;
5
6     char* stats = NULL;
7     error = neat_get_stats(mgr, &stats);
8     if (error != NEAT_OK) {
9         printf("NEAT ERROR: %i\n", (int)error);
10        return;
11    } else if (stats != NULL) {
12        printf("json %s\n", stats);
13    }
14
15    free(stats);
16 }
```

Listing 20: Simple application code that prints NEAT Statistics.

An application can leverage the information provided by the NEAT Flow Endpoint Statistics to verify that the provided Transport Services are consistent with the requested features/properties, and also to trace decisions made by the NEAT System throughout the progress of NEAT Flows which are normally not intended to be reported back to the application. For example, an application can periodically request current state information in order to be aware of any handover decisions (in case seamless handover is enabled) or changes in transport protocol parameters made by the NEAT System.

Furthermore, the statistical information provided by this building block, combined with the other types of statistics (e.g., path and interface statistics) that are also exposed through the Diagnostics and Statistics Interface, can allow applications to monitor the actual application and network performance

in order to implement more specialised functionalities. For example, this information may inform application decisions on controlling the behaviour of other applications (e.g., controlling interactions with an SDN controller) in order to optimise performance.

**Provided Transport Service Feature(s):** This building block is meant mainly for diagnostic purposes and therefore does not relate to any specific Transport Service Features.

**Related building blocks:**

- NEAT Logic (§ 3.1.3).
- Policy Manager (via Policy Interface) (§ 3.4.1).

## 3.2 NEAT Transport Components

The NEAT User API offers applications seamless access to the Transport Service Features of standardised transport protocols, while ensuring packets get through the network in the presence of non-supportive middleboxes, or challenging network paths, and providing a path to seamlessly introduce new transport protocols or transport protocol features (e.g., a Less-Than-Best Effort (LBE) Transport Service that considers data-delivery deadlines, developed in WP3 [23,25]). While the *selection* of transport protocols are handled by the NEAT Selection Components, the NEAT Transport Components are responsible for *configuring and managing* the transport protocols.

NEAT Transport components offer a consistent abstraction over underlying transport services, e.g., implementing virtual accept for UDP, thus providing a single API to applications. NEAT Transport components allow applications to use transport stacks implemented natively in operating systems and stacks implemented in userspace, such as the `usrctp` stack (§ 3.2.1).

The NEAT Transport components manage the combination of protocol parameters used (e.g., use of Nagle and other socket options to create a low-delay TCP service) and the transport protocol components that are enabled (e.g., activation of SCTP-PR or SCTP-PF respectively to create a partial reliability service or a service supporting fast path failover).

In addition, NEAT Transport components implement the necessary functionalities to assign different priorities to each NEAT Flow within a flow group. These priorities then could be used by different mechanisms at the transport level (§ 3.2.3)—e.g., coupled congestion control in TCP and per-stream scheduling priorities in SCTP.

### 3.2.1 NEAT-integrated SCTP user-space stack

The `usrctp` SCTP userland stack [11] presents a possibility to use SCTP as transport protocol when the OS does not support kernel SCTP (like Windows, NetBSD or Mac OS X), and also with Linux when SCTP is not loaded explicitly.

`Usrctp` has been ported from the FreeBSD SCTP kernel sources and is still kept in sync. A fully featured API has been provided to easily port applications.

As SCTP is an attractive alternative to TCP and UDP, it should be provided for all platforms that NEAT runs on. Therefore, we forked the original `usrctp` repository and adapted it to the needs of the NEAT System.

**Thread Support:** When using `usrctp` the transport is performed by raw sockets whose I/O is handled by threads. NEAT uses `libuv` as support library for asynchronous I/O. In order to also use this library for `usrctp` the handling of one thread per socket had to be substituted by a single-threaded algorithm. An option was introduced to `usrctp` to allow the user to either use thread support or not. In the latter case new API functions were also necessary to open the sockets from “outside” `usrctp`. An IPv4 and an IPv6 raw socket are opened from NEAT at startup and their handling is handed over to `libuv`. A new timer is also managed by `libuv` that polls the SCTP event loop.

Whenever data arrive from the network, one of the sockets becomes readable. When the timer expires, `libuv` is informed and triggers a callback function in NEAT that calls a function in `usrctp` to receive the data waiting at the socket. In the opposite direction, data is sent directly by writing the complete message (including the IP header) to the socket.

**UDP encapsulation:** SCTP is not as widely deployed as TCP and UDP, therefore it can sometimes be blocked by middleboxes. To enable SCTP transport nevertheless, SCTP packets can be encapsulated in a UDP datagram and sent as UDP payload to the receiver.

To enable this feature in NEAT two more sockets (UDP/IPv4 and UDP/IPv6) have to be opened, which are handled in the same way as the raw sockets. One advantage of UDP encapsulation is that no root privileges are needed which is the case when raw sockets are used. That is the reason why often SCTP/UDP is preferred over pure SCTP.

The user can choose to send data over SCTP/UDP instead of SCTP by choosing SCTP/UDP as transport property. NEAT then opens a `usrctp` socket and sets the socket option `SCTP_REMOTE_UDP_ENCAPS_PORT` to specify the UDP encapsulation port. Afterwards the connection is handled like a normal SCTP association.

**Handling Upcalls:** In addition to the sockets that handle the traffic between `usrctp` and the network, a `usrctp` socket is needed to manage the data transfer between NEAT and `usrctp`. This interface is defined by the `usrctp` socket API.

In the NEAT framework the sockets between NEAT and the transport layer are managed by `libuv`, i.e., `libuv` polls for a socket status change and calls a callback function whenever the socket becomes readable, writable or when an error occurs. `Libuv` expects sockets to be represented by file descriptors, however, the opening of a `usrctp` socket returns a pointer to a struct. Therefore, `libuv` cannot be used to handle `usrctp` sockets. The communication between `usrctp` and its upper layer is either handled by callback functions or directly by using `send` and `receive` calls whenever it seems adequate. In accordance with the `libuv` callback functionality we implemented a new upcall handling. An upcall function can be registered that is called whenever the socket becomes readable, writable or when an error occurs. To ensure stability the function is only triggered at the end of the handling of an incoming packet. Now, it is up to the user to call the `usrctp` `send` and `receive` functions. The callback function has to handle all events concerning I/O processes and call the corresponding operations that were registered by the user.

**Multihoming for SCTP:** An outstanding feature of SCTP is the use of multiple paths for one association, called multihoming. This means that during one connection separate paths between the peers can be used either as backup in case of a path failure or for sharing the load between them.

A first step to support multihoming in NEAT for SCTP and `usrctp` is to bind several addresses to one socket and thus reduce the number of candidates that have to be tested by Happy Eyeballs.

The candidates created by Happy Eyeballs are scanned and only those are taken where the source address corresponds to one of the addresses the user wants to send from. Then for each destination address all source addresses are assembled and bound with the `bindx` command which reduces the number of candidates to the number of destination addresses.

To enable the user to use multihoming and set the source addresses to be bound, two properties have been introduced, one to enable multihoming and one to set the addresses as a comma-separated list. The destination addresses can also be listed like this.

**ICMPv4 and ICMPv6 Support:** Using Happy Eyeballs implies that many potential connections might be tested for several protocols on different platforms. To speed up this process it is necessary to get the information immediately, whether the destination is available for a specific protocol or port or whether it is reachable.

To enable interpretation of ICMP messages that are sent following a failure by the destination, the ICMP messages must be handed to the transport protocol. Special raw sockets for ICMPv4 and ICMPv6 have to be opened that can receive ICMP datagrams. These messages are forwarded to SCTP to be interpreted so that the adequate measures can be taken.

Support for ICMPv4 and ICMPv6 has been integrated in `usrctp`, so that the upper layer is notified when a connection is aborted because of an ICMP *destination unreachable* message. For the integration in NEAT, two additional sockets have to be opened and handled by `libuv`, and the upcoming connection failure indication has to be interpreted correctly.

**Examples of the operation:** All existing examples in NEAT repository will use `usrctp` instead of kernel SCTP, if the `USRCTP_SUPPORT` option is set for NEAT. UDP encapsulation can be chosen by selecting SCTP/UDP as transport property. A new example program has been added to illustrate the use of multihoming, i.e., setting the necessary properties. The extended userland stack is available via <https://github.com/NEAT-project/usrctp-neat>.

**Provided Transport Service Feature(s):** This building block is part of the most basic functionality of the NEAT System and it does not relate to any specific Transport Service Feature.

#### Related building blocks:

- NEAT Logic (§ 3.1.3).

### 3.2.2 Middlebox Traversal

To allow NEAT-based applications to operate also in a peer-to-peer scenario, middlebox traversal has to be considered. Modern Web Browsers support WebRTC, which allows a peer-to-peer communication supporting audio, video and data channels. The data channels provide a flexible message-oriented communication, allowing the user to control the reliability and message sequence preservation. DTLS is used to provide security. WebRTC uses STUN [36] to determine the set of usable addresses, ICE [35] to ensure reachability and TURN [32] in case of very restrictive communication possibilities. These are the state of the art protocols for implementing middlebox traversal.

Therefore, it seems appropriate to integrate the data channel protocol stack used by WebRTC into NEAT. This provides middlebox traversal capabilities and the flexibility of the data channels can be used to provide the required services for NEAT flows. Using the WebRTC data channel stack within

NEAT should also allow NEAT applications to communicate with applications running in a modern Web Browser using WebRTC.

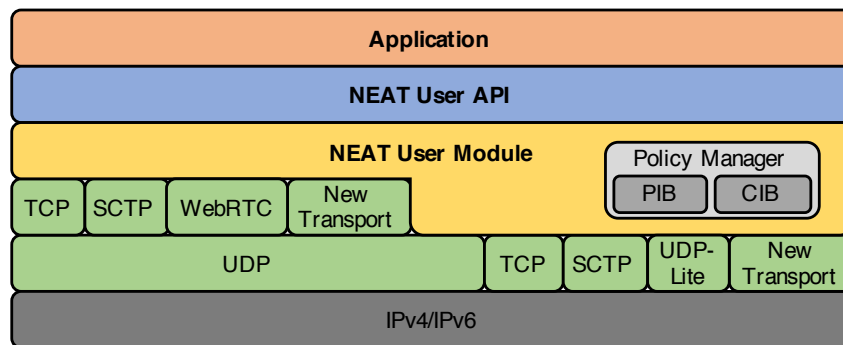


Figure 7: The architecture of the NEAT System including WebRTC.

Figure 7 shows the NEAT architecture as depicted in [38]. WebRTC performs as a new transport protocol over UDP. However, it cannot be looked upon as a single protocol like TCP or SCTP; several other protocols are needed to make WebRTC work.

**RawRTC Library:** To be able to integrate WebRTC in NEAT, a library was needed that fulfills the task to set up a secure connection with the peer.

We chose the RawRTC library [7] that adds SCTP, DTLS and data channel support to libre [2], a portable and generic library for real-time communications with asynchronous I/O support. The RawRTC library ships with its own version of the usrsctp library.

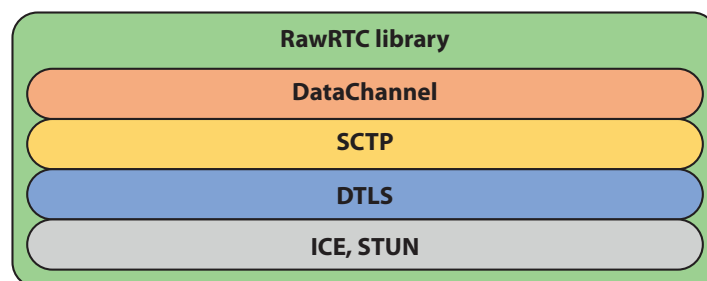


Figure 8: The RawRTC library

Figure 8 shows the four layers of the RawRTC library that are needed to realize the functionality of WebRTC.

**Integration of RawRTC in NEAT:** When integrating another transport protocol into the NEAT framework, two interfaces have to be adapted.

First, the library had to be connected to the lower layer, i.e., to UDP. The RawRTC library relies on libre's event loop for its I/O. As NEAT's communication is based on the libuv library, the event loop of RawRTC had to be substituted by the libuv event loop. This means that all sockets had to be changed to be managed by libuv. This included also the timers that trigger the polling. To achieve this, functionality in libre had to be changed.

Second, the NEAT framework had to be connected to RawRTC. The communication to set up a peer connection includes the gathering of ICE candidates, their exchange with the remote peer and

the starting of a secure transport connection. As soon as an SCTP association is established, a data channel can be opened.

As ICE gathers the candidates for the connection and sets it up, the NEAT integrated algorithm for Happy Eyeballs is not needed. Therefore, a “shortcut” is taken to bypass it.

Also the `neat_read` and `neat_write` functions are not needed to actually read from and write to the socket. The data from the user is delivered to RawRTC and transferred to the data channel which on the way back receives the user data.

To handle this functionality a shim layer had to be implemented to provide the compatibility between the NEAT framework and the RawRTC library.

**API Considerations:** The means of transport in WebRTC are data channels. To enable the user to have influence on the number of data channels and to be able to open and close them, they were mapped on flows. This means that whenever the user opens a new flow, a data channel is created, which results in the opening of a new channel on the peer’s side.

Our aim was to change the API as little as possible. Therefore, the port number in the `neat_open` call is used to distinguish between the caller and the callee by setting it to 0 or 1, respectively. A “listening” flow is created for both peers to function as a master for the other flows and, thus, for the data channels.

An example application has been added in the public Github repository to show the approach<sup>10</sup>.

**WebRTC Signaling Server:** Establishing a WebRTC Peer-Connection between two peers requires a third-party signaling channel to resolve how to establish the connection over the network. To establish a WebRTC Peer-Connection, both peers have to exchange their address candidates and their connection parameters via SDP [24]. WebRTC does not specify a transport mechanism for the signaling information. Since NEAT follows the approach to make the usage of transport protocols as easy as possible for the application developers, it ships with an integrated signaling client, which may be used. As an alternative to the integrated signaling client, application developers can integrate their own signaling mechanism.

The integrated signaling client communicates via TCP with NEAT supporting signaling server, it supports native NEAT clients as well as browser-based WebRTC clients via socket.io [9]. The signaling uses the popular node.js [4] JavaScript framework and supports multiple platforms, including FreeBSD, Linux and macOS.

**Provided Transport Service Feature(s):** There are no specific Transport Service Features associated to this building block.

**Related building blocks:**

- NEAT Logic (§ 3.1.3).
- Happy Eyeballs (§ 3.3.1).

---

<sup>10</sup>[https://github.com/NEAT-project/neat/blob/neat-rawrtc/examples/peer\\_webrtc.c](https://github.com/NEAT-project/neat/blob/neat-rawrtc/examples/peer_webrtc.c)

### 3.2.3 Local flow priority

The NEAT System provides a way to prioritise the data transfer between NEAT Flows. This is performed by assigning different priorities to the flows belonging to a certain *flow group*.

The rationale behind using flow groups is to couple flows that might share a bottleneck and control their sending rate (either by scheduling or by setting the transport's congestion window) in such a way that they can achieve their share of bandwidth without necessarily competing with each other over the available capacity. The decision to choose a certain flow group is taken at the API level and by the user (e.g., for flows that share the same source and destination).

The flow group is an integer-type number that can be assigned to a certain NEAT flow and is specified by the user with a `NEAT_TAG_FLOW_GROUP` optional argument `tag`. NEAT flows with the same flow group will be coupled together and their respective data will be transferred according to their relative priority within the group, based on the underlying mechanism used at the transport protocol level. This priority is defined by a float-type `NEAT_TAG_PRIORITY` optional argument `tag`. The optional arguments are set using the `NEAT_OPTARG_INT` and `NEAT_OPTARG_FLOAT` macros respectively and are passed to `neat_open` upon opening a NEAT flow.

Multiple flows can be created and be assigned the same flow group number. Let us assume that the kernel supports the Coupled Congestion Control (CCC) mechanism [26], and that Happy Eyeballs returns a TCP connection for each flow. In this case, the priorities are used to couple the congestion controllers of NEAT flows (i.e., TCP connections) belonging to the same flow group by giving them sending rates (i.e., congestion windows, *cwnds*) proportional to their priorities. If the kernel does not support the CCC mechanism, or if Happy Eyeballs returns SCTP handles for some flows but TCP for others, the priorities are ignored by NEAT and the *cwnd* values are determined by each flow's congestion control mechanism individually. The CCC mechanism is already implemented for the FreeBSD kernel<sup>11</sup>.

When the transparent flow mapping feature (developed in WP3 and described in D3.2 [23]) is supported and the value returned by `happy-eyeballing` uses SCTP, each NEAT flow is transparently mapped into a SCTP stream within a SCTP association. As a result all NEAT flows using the same SCTP association between the same end-points would use a common congestion controller. The priorities are therefore translated into per-stream SCTP scheduling priorities as defined in [37]. If supported by the SCTP stack, a weighted fair queueing scheduler (WFQ) between the SCTP streams is used. The weight is configurable per outgoing SCTP stream. This scheduler considers the lengths of the messages of each stream and schedules them in a specific way to use the capacity according to the given weights.

Deliverable D3.2 [23] elaborates on the underlying CCC mechanisms used for flow-group prioritised transfer for TCP and SCTP.

**Some examples of the operation:** Assume a scenario where there are two flows, one with `NEAT_TAG_PRIORITY` of 0.75 starting at  $t = 3$  s, and a second flow with `NEAT_TAG_PRIORITY` of 0.25 starting at  $t = 30$  s, with both being assigned the same flow group `NEAT_TAG_FLOW_GROUP` as 1.

To test this, we adapted the `tneat`<sup>12</sup> traffic-generation tool to accept flow group priority and number as command-line arguments as shown in Listing 21 (`-w` for priority and `-g` for flow group). Each

<sup>11</sup>Kernel support code is available at <https://naeemk@bitbucket.org/naeemk/freebsd11-ccc.git>.

<sup>12</sup>`Tneat` is presented in Deliverable D4.2 [14]. Code for the adapted version of `tneat` currently resides in `neat/examples/tneat_fg.c` in the `oystedal/dscp` branch of the public NEAT code repository, at <https://github.com/NEAT-project/neat>.



`tneat_fg` instance corresponds to a single NEAT flow and generates traffic by sending a certain number of messages to a specified host (`-n` and `-l` specify the number and length of messages respectively).

```
1 #! /usr/bin/env sh
2 COUPLED_FLOWS=1
3
4 sleep 3
5 if [ $COUPLED_FLOWS == 1 ]; then
6     ./neat/examples/tneat_fg -p 5001 -g 1 -w 0.75 -l 1000 -n 100000 10.0.0.5 &
7     sleep 30
8     ./neat/examples/tneat_fg -p 5002 -g 1 -w 0.25 -l 1000 -n 15000 10.0.0.5 &
9 else
10    ./neat/examples/tneat_fg -p 5001 -g 0 -w 1.0 -l 1000 -n 115000 10.0.0.5 &
11 fi
12
13 sleep 120
14 echo "$0 DONE"
```

Listing 21: Shell script using `tneat_fg` command-line with different priorities in a flow group.

Listing 22 shows how the flow group and priorities as well as the CCC algorithm (`newreno_afse`) are passed to the NEAT System as optional arguments via `neat_open`, in `tneat_fg.c`.

```
1 ...
2 static float config_priority = 1.0f;
3 static uint16_t config_group = 0;
4 ...
5 int
6 main(int argc, char *argv[])
7 {
8     ...
9     NEAT_OPTARGS_DECLARE(2);
10    NEAT_OPTARGS_INIT();
11    ...
12    while ((arg = getopt(argc, argv, "l:n:p:w:g:P:R:T:v:")) != -1) {
13        switch(arg) {
14            ...
15            case 'g':
16                config_group = atoi(optarg);
17                break;
18            case 'w':
19                config_priority = atof(optarg);
20                break;
21            ...
22        }
23    }
24    ...
25    if (config_active) {
26        if (config_group) {
27            NEAT_OPTARG_INT(NEAT_TAG_FLOW_GROUP, config_group);
```

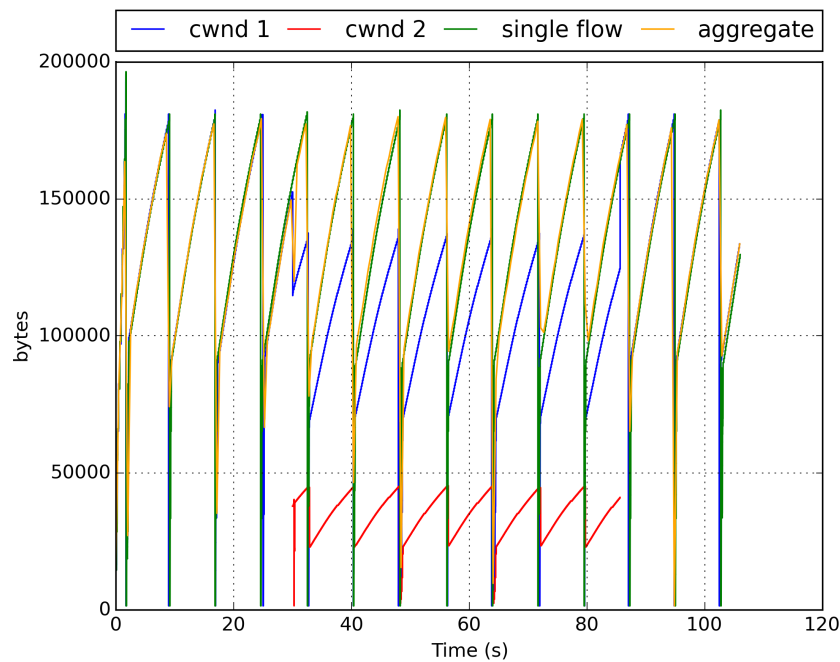


Figure 9: The *cwnd* (in bytes) plot of two TCP flows using coupled congestion control with priorities compared to a single TCP flow scenario. The *aggregate* line depicts the sum of cwnds in the two-flow scenario.

```

28     NEAT_OPTARG_FLOAT(NEAT_TAG_PRIORITY, config_priority);
29     NEAT_OPTARG_STRING(NEAT_TAG_CC_ALGORITHM, "newreno_afse");
30 }
31 if (neat_open(ctx, flow, argv[optind], config_port, NEAT_OPTARGS,
32     NEAT_OPTARGS_COUNT) == NEAT_OK) {
33     neat_start_event_loop(ctx, NEAT_RUN_DEFAULT);
34 } else {
35     fprintf(stderr, "%s - neat_open failed\n", __func__);
36     result = EXIT_FAILURE;
37     goto cleanup;
38 }
39 ...
40 cleanup:
41     if (ctx != NULL) {
42         neat_free_ctx(ctx);
43     }
44     exit(result);
45 }

```

Listing 22: Assigning flow group, priority and CCC selection in `tneat_fg.c`.

Figure 9 presents the results obtained with the example in Listing 21, using a real-life test on a 10 Mbps bottleneck link with RTT=100 ms (typical for access link scenarios), with a FreeBSD machine acting as the sender and using NewReno as the underlying TCP congestion control mechanism with coupled congestion control enabled.

It can be seen in Figure 9 that between  $t = 30$  s and  $t = 85$  s (the duration when two flows coexist),

Table 1: Transport protocols and NEAT security in the Core Transport System prototype. *All* means all the OSs supported by NEAT at the time of writing, as listed in § 2.1.1.

Security stack	Implemented in	OS platform
TLS/TCP/IP	kernel	all
DTLS/UDP/IP	kernel	all
DTLS/SCTP/IP	kernel	Linux and FreeBSD
DTLS/SCTP/IP	user-space	unsupported
DTLS/SCTP/UDP/IP	kernel	FreeBSD
DTLS/SCTP/UDP/IP	user-space	unsupported
SCTP/DTLS/UDP/IP (with WebRTC, see § 3.2.2)	user-space	all

flow 1 gets 3/4 of the cwnd while flow 2 gets 1/4 of the cwnd, on the average. This shows that the two flows are coupled, because the coupling causes each to increase and decrease their cwnds at (close to) the same time. Outside of this interval—when flow 2 either has not started yet or has terminated—flow 1 gets all of the cwnd.

**Provided Transport Service Feature(s):**

- NEAT flow group.
- NEAT flow priority.

**Related building blocks:**

- NEAT Logic (§ 3.1.3).

**3.2.4 Security**

The Security component of the NEAT System offers end-to-end transport security to NEAT applications, including encryption, integrity, and authentication. Applications that, for reasons of backwards compatibility, cannot require full security may still utilize opportunistic security though this is not yet available.

The NEAT System provides secure connections using the TLS [17] and DTLS [34] protocols. TLS is used over TCP and DTLS for SCTP and UDP. As shown in Table 1, the possible combinations that NEAT can provide are: TLS/TCP/IP, DTLS/UDP/IP, and DTLS/SCTP/UDP/IP or DTLS/SCTP/IP for a kernel-level SCTP stack. The OpenSSL library [6] that offers TLS and DTLS support is used in the current implementation. NEAT security components support TLS/TCP/IP and DTLS/UDP/IP on all supported OS platforms, as well as DTLS/SCTP/IP for a kernel-level SCTP stack on Linux and FreeBSD. DTLS/SCTP/UDP/IP is supported on FreeBSD as well when the SCTP/UDP/IP combination is provided by the FreeBSD kernel. DTLS/SCTP/IP or DTLS/SCTP/UDP/IP for a user-space SCTP stack requires properties currently not supported by OpenSSL, i.e., the BIO memory access interface for OpenSSL [5] does not provide a message preserving memory access which is required by SCTP, furthermore OpenSSL provides an interface that requires a kernel socket.

The use of the Security component depends on application requirements and configured policies. Transport security can be configured in one of the following ways:

- A secure connection is requested including a certificate verification.
- A secure connection is requested without a certificate verification: a certificate verification will not be performed.
- A secure connection is optional: if a secure connection cannot be established the NEAT System will not return an error, instead it will establish a new insecure connection. This will add an additional delay. This is currently not available but it is supported by the design.
- A non-secure connection is requested: this option will not involve the Security component.

This corresponds to two NEAT Flow properties:

- *security*: boolean indicating whether or not to use a secure connection.
- *verification*: must be verified, verification is optional, or do not verify. This property is only relevant if a secure connection is used and defaults to true.

A list of trusted Certification Authorities (CA) is currently hardcoded using the Mozilla CA root program data<sup>13</sup>. Future work should enable a trusted list to be specified as a policy. The server certificates and private key files are set using `neat_secure_identity` and the function accepts only the PEM format [12, 27, 28, 31]. SSL3 and RC4 ciphers are not allowed because they do not meet current best practices for communications security. TLS and DTLS can be limited to advertise and accept only certain TLS/DTLS versions and cipher suites—e.g., advertise only TLS 1.2 and the `TLS_ECDHE_RSA_WITH_AES_128_GCM_SHA256` cipher suite. The TLS/DTLS versions and cipher suites to be accepted will be defined as policies. To extend the flexibility of the NEAT System, an application may further limit or extend these lists for each individual request by setting the corresponding NEAT Flow property. This is currently not implemented.

Security is added in the NEAT as an additional protocol above the operating system socket. The NEAT Logic calls into the Security building block to perform data encryption and decryption before data is written to/read from an operating system socket. During a connection establishment the Security building block performs TLS or DTLS handshake and certificate verification depending on application requirements.

The NEAT System can be extended to use TCPINC [10] although this is not part of the project work.

**Some examples of the operation:** The security component is used automatically by the NEAT Logic if a secure connection is required. An application needs to specify the corresponding property. Listing 23 shows a simple example of how to request a secure connection to a peer. This example is based on Listing 17; an application only needs to change flow properties to request a secure connection.

```
1 int main(int argc, char *argv[])
2 {
3     if ((ctx = neat_init_ctx()) == NULL) {
4         debug_error("could not initialise context");
5         result = EXIT_FAILURE;
6         goto cleanup;
7     }
8 }
```

<sup>13</sup><https://www.mozilla.org/en-US/about/governance/policies/security-group/certs/>

```
9 // new neat flow
10 if ((flow = neat_new_flow(ctx)) == NULL) {
11     debug_error("neat_new_flow");
12     result = EXIT_FAILURE;
13     goto cleanup;
14 }
15
16 char *config_property = "{\
17     \"transport\": [\
18         {\
19             \"value\": \"SCTP\", \
20             \"precedence\": 1\
21         }, \
22         {\
23             \"value\": \"SCTP/UDP\", \
24             \"precedence\": 1\
25         }, \
26         {\
27             \"value\": \"TCP\", \
28             \"precedence\": 1\
29         } \
30     ] \
31     \"security\": {\
32         \"value\": true, \
33         \"precedence\": 2 \
34     } \
35 }";
36
37 // set properties
38 if (neat_set_property(ctx, flow, config_property)) {
39     fprintf(stderr, "%s - error: neat_set_property\n", __func__);
40     result = EXIT_FAILURE;
41     goto cleanup;
42 }
43 ...
```

Listing 23: Use of the NEAT Security component.

**Provided Transport Service Feature(s):**

- NEAT flow security.

**Related building blocks:**

- NEAT Logic (§ 3.1.3).

### 3.3 NEAT Selection Components

The NEAT Selection components provide functions that map the requirements provided by the application to one or more transport endpoints and a set of transport components that can realise the

---

**Algorithm 1** NEAT Happy Eyeballs Algorithm

---

```
1: procedure HAPPY_EYEBALLS_COMPONENT(in listOfCandidates : list of transport solutions)
2:   Require: listOfCandidates is sorted in priority order and len(listOfCandidates) > 0
3:   currentCandidate ← listOfCandidates.first()
4:   repeat
5:     if getPriority(currentCandidate) > 0 then
6:       delta ← convertToTimeInterval(getPriority(currentCandidate))
7:       scheduleAt(now() + delta, doAsynchConnectionAttempt(currentCandidate, connectionCallback))
8:     else
9:       doAsynchConnectionAttempt(currentCandidate, connectionCallback)
10:    end if
11:    currentCandidate ← listOfCandidates.nextCandidate(currentCandidate)
12:  until currentCandidate = endOfList(listOfCandidates)
13: end procedure

14: procedure CONNECTION_CALLBACK(in candidate : transport solution, out connection : transport connection)
15:   if connection ≠ NONE then
16:     policyManager.cacheResultConnectionAttempt(candidate, SUCCESS)
17:   else
18:     policyManager.cacheResultConnectionAttempt(candidate, FAILURE)
19:   end if
20: end procedure
```

---

required Transport Service. These functions are provided now by the Happy Eyeballs (§ 3.3.1) and Happy Apps (§ 3.3.2) components. While Happy Eyeballs provides the necessary transport selection functionalities at the transport layer using the information provided by the transport, Happy Apps provides a means for selection of appropriate protocol and configurations using feedback from the application when such feedback is unavailable at the transport layer level.

### 3.3.1 Happy Eyeballs

The Happy Eyeballs (HE) building block is part of the NEAT User Module and comprises one of the NEAT Selection components. A description of the main parts of this building block is the subject of an IETF Internet Draft, authored by project participants [21].

In the first part of the NEAT selection process (i.e., steps 2–8 in Figure 4), the Policy Manager combines requirements from an application obtained through the NEAT User API, with available transport protocols, transport-protocol parameters, and feasible transport endpoints, i.e., IP addresses and port numbers. Together, they are used to create a list of candidate *transport solutions*. Each transport solution on the list has a priority. The priority is a positive integer with zero being the highest priority. The list is sorted in ascending order on the basis of the priority.

The pseudo-code for the Happy Eyeballs building block is presented in Algorithm 1. The component takes as input the list of candidate transport solutions created in the first part of the NEAT selection process. Next, it traverses the list and schedules for the transport solutions to be tried out after a certain amount of time  $t$  has elapsed, governed by their priority  $p$ , as:  $t = p \times \Delta$ , where  $\Delta$  is a fixed time interval which has been heuristically determined. The transport solution which is able to establish a connection first, is the one selected by the Happy Eyeballs building block, and the one whose handle is returned.

As remarked in RFC 6555 [42], a Happy Eyeballs algorithm should not waste networking resources by routinely making simultaneous connection attempts. To this end, the Happy Eyeballs component instructs the Policy Manager to *cache* the outcome of previous connection attempts. Cached connection attempts are valid for a pre-set time after which they become invalid and have to be repeated. The caching lifetime is at present hardcoded to a fixed value, but the library will be updated so that

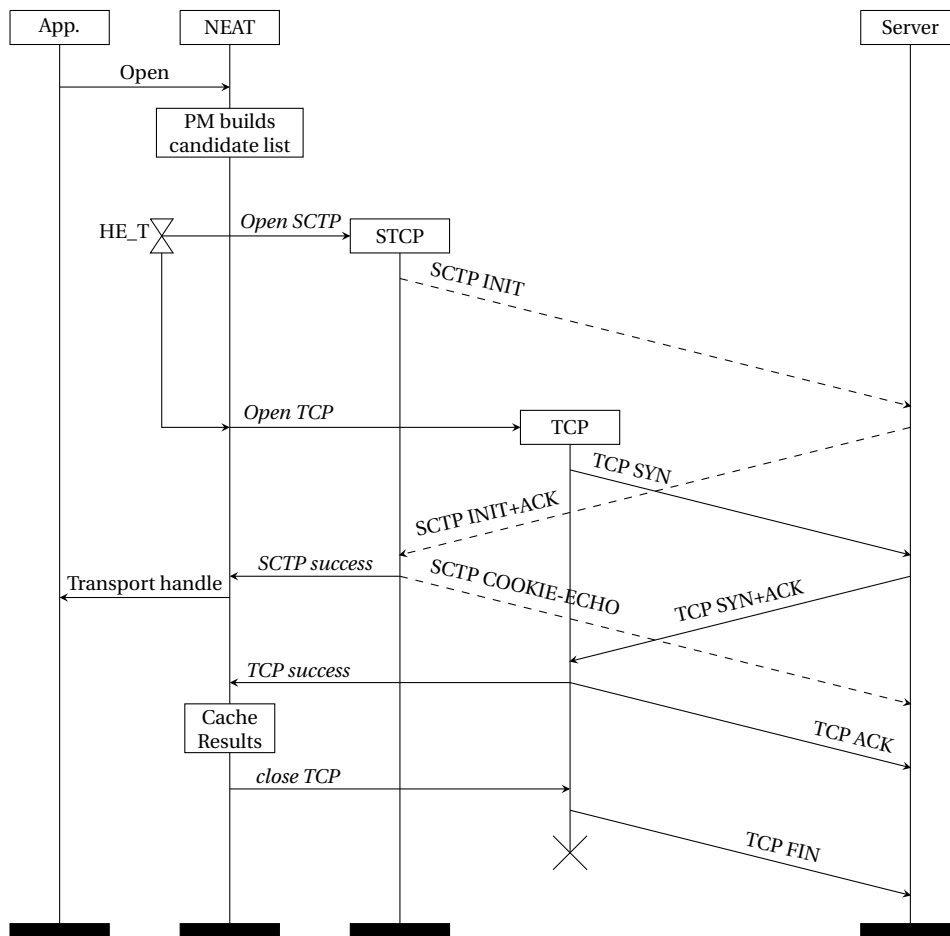


Figure 10: Message sequence chart illustrating the NEAT Happy Eyeballs transport selection process when selecting between TCP and SCTP, with SCTP preferred.

this lifetime is part of the system configuration. Our experimental work [33] suggests a significant reduction in terms of CPU load with caching. We observed that the CPU load decreases linearly with increasing cache hit-rate, and resulted in a more than 40% reduction of CPU load for unencrypted traffic and almost a 20% reduction for encrypted traffic. In terms of memory, our work reported in [33] concluded that Happy Eyeballs only has a marginal impact on kernel memory usage.

**Some examples of the operation:** As an example of how Happy Eyeballs works, consider the scenario illustrated in Figure 10. Both the client and server support TCP and SCTP. The Policy Manager puts together a list of candidate transport solutions. The NEAT Logic calls Happy Eyeballs, which traverses the candidate list, and makes asynchronous connection attempts for each candidate transport solution. Since the TCP transport solution has a lower priority than the SCTP transport solution, the connection attempt for the TCP transport solution is delayed with respect to the one for SCTP. As detailed before, the length of the delay depends on the priority. In the callback routine invoked when a connection attempt returns, the outcome of the connection attempt (success or failure) is cached by the Policy Manager through the Policy Interface. Since both connection attempts succeed, they

are cached as successful connection attempts. In the example depicted in Figure 10, the SCTP connection attempt completes before the TCP connection attempt, and Happy Eyeballs returns the SCTP connection.

**Provided Transport Service Feature(s):**

- NEAT selected transport protocol.

**Related building blocks:**

- NEAT Logic (§ 3.1.3).
- NEAT Policy Manager (via Policy Interface) (§ 3.4.1).
- Connect to a name (§ 3.1.4).

**3.3.2 Happy Apps (application-level feedback mechanisms)**

Happy Apps (HA) is part of the NEAT User Module and offers selection mechanisms when the underlying transport protocol does not provide the signals required by the NEAT Logic. Non-connected network transports such as UDP and UDP-Lite do not provide mechanisms via the socket API to signal if traffic is making it to the end host, instead it is up to the application-level protocol to decide if it is “happy” with the current flows provided by the NEAT System. The application is the only entity in a datagram transport using UDP(-Lite) that knows the reception state of the receiver, and hence whether a path is operating as required.

The Happy Apps mechanisms allow applications to take advantage of network properties that might not always work while still using the single NEAT User API.

Happy Apps can support fallback queries for different types of network issue, specified to NEAT with policy requests. Currently *QoS Fallback*, further explained below, is specified and implemented as a fallback mechanism in the NEAT System, but there are many possible Fallback Services an application could subscribe to. Some examples are:

- **Abstract QoS Fallback with DiffServ:** The NEAT System offers a mechanism that allows the application to signal its QoS requirements and match these to expected network support for QoS. This can allow NEAT policy to be utilised to select the most appropriate DSCP marking (or other QoS service). The mapping can take advantage of local configuration data, and QoS information provided by the network. By providing fallback through potential candidates it can help ensure that application traffic makes its way through the network with the most appropriate QoS available. This can be useful when a specific feature is blocked by some network boxes on the path.

Abstracting selection to the stack also allows the stack to utilise other path information that it gathers (via a CIB source utilising the Policy Manager) to help make this selection, and enables evolution of the selection process independent of the application code.

- **Custom SDN Provisioning Requests:** An application running in an SDN environment requires a large amount of bandwidth to perform a backup operation. Depending on network load, either the standard interface is used or a special link can be allocated. The application manages to maintain the required rate at the start of the session, but the bandwidth falls off. When the application signals it is “unhappy” to the NEAT System, the SDN Controller dynamically allocates an assured link for the application to use.



- **Firewall Traversal:** An application requires a service that is operated behind a firewall that only accepts admitted traffic. The application uses NEAT to connect to the firewall, but the traffic is stopped, NEAT falls back to a protocol that supports negotiation with the firewall and authenticates the application allowing traffic through.

The Happy Apps system is intended to be run throughout the application's lifetime; path changes that break certain network parameters can occur at any time. The NEAT System improves the API by giving applications a way to keep flows alive with mostly automatic fallback.

It is important to note that the Happy Apps mechanism may not be able to make any further changes to network parameters on behalf of the NEAT-based application, e.g., the application falls through all the selected QoS values without finding one that satisfies its requirements. In this case HA will remain in the final selected state until the application either decides to timeout and shutdown or it attempts to vary parameters itself.

**Abstract QoS:** NEAT provides mechanisms for transports to use QoS within the network. QoS is implemented in NEAT by setting the DSCP bits in the IP TOS field. Applications built with the socket API can access this field via a `setsockopt`, NEAT exposes the QoS field directly in the NEAT Flow and handles parameter setting in the socket stack on behalf of the application.

Applications using NEAT can set a DSCP value directly or they can use the *Abstract QoS* facilities offered by the NEAT System. NEAT offers abstract QoS types to give application developers access to QoS with a higher level API.

The abstract QoS values in NEAT allow application developers to request a QoS service which is comparable to the requirements of the application they are building. Rather than an application developer directly requesting a code point such as EF (Expedited Forwarding), with the NEAT System an application developer can request a high level `NEAT_QOS_AUDIO_H1` (a class representing high bandwidth audio traffic) and the NEAT System will perform the mapping down to the concrete QoS.

NEAT Abstract QoS uses the Policy Manager when mapping between the Abstract QoS value and the concrete DSCP mark to be used for the packet. Selection with the Policy Manager allows the NEAT System to offer dynamic mapping from abstract to concrete QoS while using information generated by other flows. Table 2 shows the default map from abstract QoS to concrete DSCP values in the NEAT System.

**Some examples of the operation:** As an example of how Happy Apps works, consider an application that requests QoS on a network path which drops all traffic marks with the AF4x DSCP. The NEAT application requests a NEAT Abstract QoS service of `NEAT_QOS_INTERACTIVE_VIDEO_H1` (Interactive Video, High Bandwidth) and sets the `feedback_query` callback. The Policy Manager evaluates the requested QoS and generates a concrete mapping to a code point, for this example AF41.

The NEAT System automatically adds the QoS fallback property, and sets up the Happy Apps mechanism because the application has requested an abstract QoS type.

NEAT Logic sets the resolved DSCP on the socket underlying the application's NEAT flow. Each datagram generated by the application via `neat_write` will generate an IP packet with the specified DSCP set. The NEAT Application can now work through its own protocol, reading and writing against the NEAT flow as desired.

The `feedback_query` callback will be triggered by the NEAT System when the feedback query interval time expires, by the default policy triggers this after around 1 second. The NEAT Application

Table 2: Possible Abstract QoS to DSCP Mappings in NEAT. Some traffic classes such as Video can have several different capacity requirement levels, the NEAT System exposes these with Very Low, Low, Medium and High capacity requirements. Applications can also request Admitted access, classes that can be guaranteed by the network with policy or dynamic provisioning.

Abstract Name	DSCP Code	DSCP Value
NEAT_QOS_AUDIO_VL	CS1	0x08
NEAT_QOS_AUDIO_L	DF	0x00
NEAT_QOS_AUDIO_M1	EF	0x2E
NEAT_QOS_AUDIO_H1	EF	0x2E
NEAT_QOS_INTERACTIVE_VIDEO_VL	CS1	0x08
NEAT_QOS_INTERACTIVE_VIDEO_L	DF	0x00
NEAT_QOS_INTERACTIVE_VIDEO_M1	AF42	0x24
NEAT_QOS_INTERACTIVE_VIDEO_M2	AF43	0x26
NEAT_QOS_INTERACTIVE_VIDEO_H1	AF41	0x22
NEAT_QOS_INTERACTIVE_VIDEO_H2	AF42	0x24
NEAT_QOS_NON_INTERACTIVE_VIDEO_VL	CS1	0x08
NEAT_QOS_NON_INTERACTIVE_VIDEO_L	DF	0x00
NEAT_QOS_NON_INTERACTIVE_VIDEO_M1	AF32	0x1C
NEAT_QOS_NON_INTERACTIVE_VIDEO_M2	AF33	0x1E
NEAT_QOS_NON_INTERACTIVE_VIDEO_H1	AF31	0x1A
NEAT_QOS_NON_INTERACTIVE_VIDEO_H2	AF32	0x1C
NEAT_QOS_DATA_VL	CS1	0x08
NEAT_QOS_DATA_L	DF	0x00
NEAT_QOS_DATA_M1	AF11	0x0A
NEAT_QOS_DATA_H1	AF21	0x12
NEAT_QOS_BROADCAST	CS3	0x18
NEAT_QOS_REALTIME_INTERACTIVE_DATA	CS4	0x20
NEAT_QOS_IMMERSIVE_AUDIO	AF41	0x22
NEAT_QOS_IMMERSIVE_VIDEO	AF41	0x22
NEAT_QOS_STREAMING	AF31	0x1A
NEAT_QOS_BACKGROUND	CS1	0x08
NEAT_QOS_ADMITTED_AUDIO	EF	0x2E
NEAT_QOS_ADMITTED_VIDEO	AF42	0x24
NEAT_QOS_ADMITTED_IMMERSIVE_AUDIO	AF42	0x24
NEAT_QOS_ADMITTED_IMMERSIVE_VIDEO	AF42	0x24
NEAT_QOS_ADMITTED_DATA	AF42	0x24

must service the callback based on its own state. Because the network drops the traffic from the NEAT Application, the application cannot say it is “happy” with the current setup and returns `NEAT_ERROR`.

NEAT Logic receives the negative application feedback, the established DSCP mark is fed into the Policy Manager so that signal can be used to avoid this for any later connections. The NEAT Logic requests another concrete DSCP mark to use for the flow, in this example the next mark to try is AF31.

The application continues to perform writes with the new QoS value, this time the mark is able to transverse the network. The `feedback_query` callback is triggered again, this time the application is seeing return data from its peer and is able to signal success by returning `NEAT_OK`.

**Provided Transport Service Feature(s):**

- NEAT flow DSCP Support.

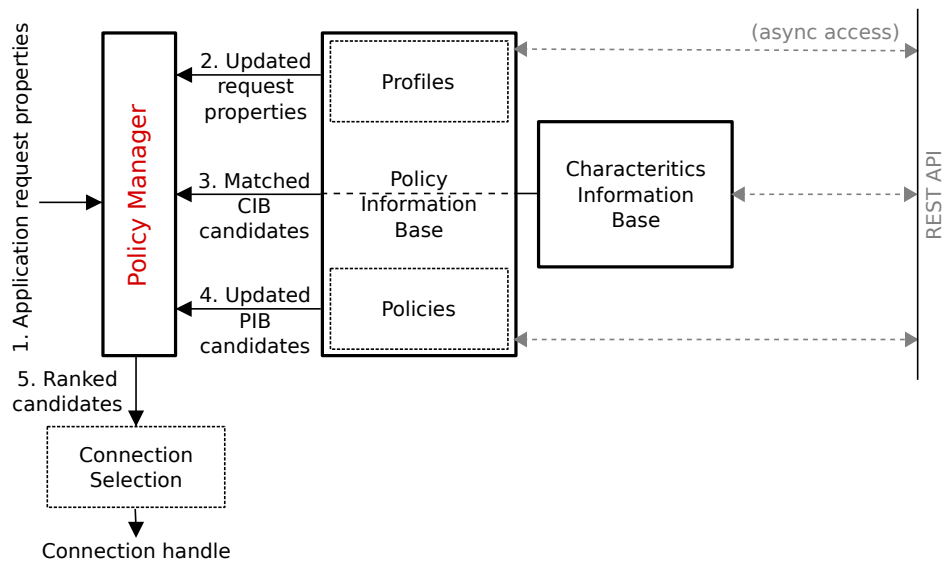


Figure 11: NEAT Policy components and their interactions.

**Related building blocks:**

- NEAT Logic (§ 3.1.3).
- NEAT Policy Manager (via Policy Interface) (§ 3.4.1).

**3.4 NEAT Policy Components**

The NEAT Policy components allow the definition of policies which shape the attributes and constraints of each new given NEAT, flow based on the properties requested by applications as well as known system and network characteristics. As a result the NEAT System can seamlessly adapt to a wide range of scenarios using administrative policies as well as information collected about the host and the attached networks.

The Policy components are comprised of the following building blocks, depicted in Fig. 11 and described in the sequel:

- Policy Manager (PM).
- Policy Information Base (PIB).
- Characteristics Information Base (CIB).

For each new connection requested by a NEAT-enabled application, the PM is responsible for generating a list of candidate transport solutions which meet the requirements defined by the application. To this end, the PM takes into account all known information about the network stored in the CIB as well as any applicable policies defined in the PIB.

The CIB acts as a repository storing information about available interfaces, supported protocols, network properties and current/previous connections between endpoints (generated from passively or actively acquired network metrics—see *CIB sources* in § 3.4.3).

The PIB acts as a repository that contains a collection of policies, where each policy consists of a set of rules linking a set of matching requirements to a set of preferred or mandatory transport characteristics. Policies can be added by the system administrator, external entities or applications, and have

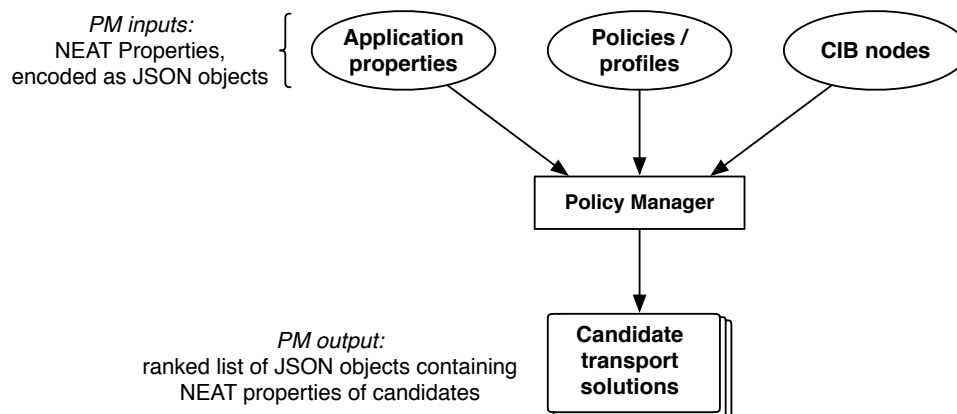


Figure 12: Policy Manager inputs and outputs.

different priorities. In addition to policies the PIB also contains so-called *profiles*, which are policies that are applied *before* the CIB lookup. Profiles are typically used to translate high-level properties into sets of concrete, low-level properties. Besides the time point of execution, policies and profiles are functionally identical.

The CIB and PIB entries are parsed to compute a list of potential candidates where each candidate contains an interface to be used, the transport protocol and associated options as well as other characteristics of the network.

The PM is responsible for implementing a strategy which (as far as possible) satisfies the given application requirements and installed policies while taking into account the knowledge about the available network resources. To achieve this the PM needs to prioritise and resolve any conflicting policies, adapting them to changing network information. The PM may optionally manage thresholds for when a CIB change triggers a new policy. Policies do not update within a NEAT Flow's lifetime, the PM is invoked only when a new NEAT Flow starts.

### 3.4.1 NEAT Policy Manager

The NEAT Policy Manager compiles all policies available at the host into a single set of valid rules. Conflicting policies are resolved based on the priority of categories: the application local policies over the global policies and external system policies over NEAT System policies. Global NEAT policies are combined during NEAT initialisation and serve as the default. Application local policies override the default policies only for the application by which these were installed. There is no lifetime associated with policies. They do not expire and have to be explicitly removed. For each application the set of valid rules compiled from policies is static and will not be changed during runtime. The PM may optionally implement some level of access control. If more NEAT components require this functionality the addition of a NEAT-wide AAA building block may be considered. The NEAT Policy Manager can expose diagnostic information about the installed and/or active policies. These statistics may be bundled with output from the NEAT Flow Endpoint Statistics (§ 3.1.5), together with statistics from CIB, to enable useful debugging.

Figure 12 summarises several concepts behind the operation of the PM, further developed below.

**PM inputs:** The PM requires three types of inputs:

- *Application properties*: an array of *NEAT properties* (see below) describing the attributes that a NEAT-enabled application desires for a newly opened NEAT Flow. The PM uses the approach described in this section (under **Some examples of the operation**, below) to extract the most suitable candidates for the property requirements. The NEAT property array is passed to the Policy Manager by the NEAT Logic through the Policy Interface.
- *Policies*: they are installed by OS developers, vendors or applications into a predefined location in the file system, or transferred through the PM REST API. Profiles are a special type of policies which are used to expand high-level properties into more concrete properties *before* the CIB lookup. Policies and Profiles use an identical syntax using the JSON format [18] (the policy format is described in § 3.4.2).
- *CIB nodes*: they provide information about transport and interface characteristics in a predefined location in the file system or transferred through the PM REST API. Each CIB node is generated and maintained by so-called *CIB sources*, which may be system local or remote (§ 3.4.3). The CIB repository is generated by all available CIB nodes which form a graph structure. Effectively the CIB expands the branches of this structure to create rows containing potential connection candidates.

**PM outputs:** The output of the PM is a ranked JSON list containing a set of candidate transport solutions and parameters for use by the NEAT Logic. Each candidate is an array of NEAT properties. Each candidate includes at least the local interface, the local and remote IP addresses and ports and the transport protocol. In addition, each candidate in the output list contains information indicating:

- Which application requirements (properties) are satisfied for the given interface/protocol/destination tuple.
- Which of the properties specified by the related evaluated policies have been verified and applied (added to the candidate).
- Which of the properties specified by the related evaluated policies the PM was not able to verify.

**NEAT Properties:** The NEAT policy system is based around the notion of NEAT properties. These are essentially `key|value` pairs describing attributes used by the components of the NEAT Policy Manager. NEAT properties are used to express attributes describing the local host and the attached networks, user requirements, or constraints imposed by policies. Examples include interface names, types, supported protocols and parameters, or topology metrics (e.g., available bandwidth).

Each property has a `key` string and a `value`. Currently property values can be:

1. A **single** boolean, integer, float or string value, e.g., `2`, `true`, or `"TCP"`.
2. A **set** of values: `[100, 200, 300, "foo"]`.
3. A numeric **range**: `{"start":1, "end":10}`.
4. **ANY** value: `null`.

Each property is further associated with a `precedence` which identifies the “importance” of the property. Specifically, the precedence indicates if the property may be modified by the Policy Manager logic or if it is immutable. Currently two property precedence levels are defined in order of decreasing priority:

**[immutable]** (precedence 2): These are mandatory properties which value cannot be changed.

**(requested)** (precedence 1): These are optional properties whose value may be overwritten. A mismatch of such properties will result in a penalty in the ranking within the PM. Such penalties are recorded as the `score` of the property.

The property's numeric `score` attribute is used to indicate the weight of the property with respect to other properties. Whenever two properties are compared and match, the score of the result is the sum of the individual scores. As a consequence the PM can use the score to determine the most suitable NEAT connection `candidate` for a given request.

In the sequel we use the following shorthand notation: property key/value pairs are separated by the `|` character and indicate the property's precedence by the bracket types shown above. We append the score to the brackets and omit it if zero. For example

- `[transport|TCP]+2`: the transport protocol *must* be "TCP"
- `(MTU|1500, 2000, 9000)`: one of the specified MTU values should be chosen *if* possible
- `(capacity|10-1000)+1`: the interface capacity should be within the numeric range specified by the integers

### **Property Operators:**

1. Comparison: `p1 == p2`

Any two NEAT properties are considered equal if their keys are identical and the intersection of their values is a non-empty set. Precedence and scores are ignored when testing for equality. A comparison — or **match** — of two properties yields a boolean result.

For instance, the operation `[transport|TCP]+1 == (transport|TCP)+3` yields true. Set and range value attributes are also considered equal if their values overlap, i.e., the expressions `[transport|TCP,UDP,MPTCP] == [transport|TCP]`, or `[latency|1-100] == [latency|55]` both return true.

2. Update: `p1 <= p2`

In the lookup workflow in the PM, properties from various sources will be compared and potentially *updated*. Updates are the mechanism used to filter candidates which violate policies or contradict information stored in the CIB. Essentially, an update is successful if the property value of the involved properties do not contradict each other.

More specifically we use the following rules. A property may only be updated by another property with the same key. A property's value may **only** be updated by another property whose precedence is greater or equal than itself — in which case it inherits the precedence of the updating property — and if both properties are not immutable (highest precedence). A property update succeeds if the ranges of the associated properties overlap — the resulting updated property will contain the intersection of the two ranges. Otherwise, if the above conditions are not satisfied, the update will fail. If the property is part of a PropertyArray a failed property update will invalidate the *entire* associated candidate.

As an example, if the immutable property `[transport|TCP]+2` is requested by an application and this property clashes with the property `[transport|UDP]` in a certain connection candidate, the candidate will be discarded.

**JSON Encoding:** The Policy Manager uses JSON to represent NEAT Properties. As an example, the property [transport|TCP, SCTP]+3 is encoded as follows:

```
"transport": {
  "value": ["TCP", "SCTP"],
  "precedence": 2,
  "score": 3
}
```

The property (latency|1-100) is encoded as:

```
"latency": {
  "value": {"start":1, "end":10}
}
```

If any property attribute is omitted, the following default values will be used:

```
"value": null, "precedence": 2, "score": 0, "evaluated": false.
```

For all externally facing interactions with the PM, NEAT properties are encoded as a JSON object. An example request containing multiple properties is encoded as:

```
1 {
2   "domain_name" : {"value": "www.google.com", "precedence": 2},
3   "port": {"value": 80, "precedence": 2},
4   "local_interface": {"value": "eth0"},
5   "local_ip": {"value": "10.10.2.14", "precedence": 2},
6   "transport": {"value": "TCP", "precedence":2, "score":2},
7   "MTU": {"value": [1500,2000,9000]},
8   "capacity": {"value": {"start":10, "end":1000}, "score":1}
9 }
```

Listing 24: JSON encoded PM query.

**NEAT PropertyArrays:** A NEAT *PropertyArray* is a simple object containing a set of NEAT properties. These objects are used to represent application *requests* received through the NEAT User API as well as the connection *candidates* generated by the PM as the request traverses the PM logic.

An example of the simplest encoding of a PropertyArray containing a set of properties looks as follows:

```
1 [{
2   "remote_ip": {"precedence": 2, "value": "10.54.1.23"},
3   "port": {"precedence": 2, "value": 8080},
4   "transport": {"value": "reliable"},
5   "MTU": {"value": [1500, 9000]},
6   "low_latency": {"precedence": 1, "value": true}
7 }]
```

To compactly express arrays with multiple property sets, PropertyArrays themselves may be comprised of several arrays containing one or more objects (MultiPropertyArrays). In this case, an `expand`

Table 3: REST API endpoints for external access to the Policy Manager.

REST resource	HTTP method	Description
/pib	GET	lists all policies installed in the host
/pib/{uid}	GET/PUT	retrieve or upload a policy with a specific Unique Identifier (UID)
/cib	GET	lists all CIB nodes installed in the host
/cib/{uid}	GET/PUT	retrieve or upload a CIB node with a specific UID
/cib/rows	GET	retrieve all rows of the CIB repository

operator yields the Cartesian product of all contained objects. For instance, a PropertyArray  $P$  containing the NEAT Properties  $p_a, p_1, p_2, p_3, p_4$  defined as  $P = [[\{p_a\}], [\{p_1, p_2\}], \{p_3, p_4\}]$  expands to two PropertyArrays  $[\{p_a, p_1, p_2\}]$  and  $[\{p_a, p_3, p_4\}]$ . Thus, MultiPropertyArrays may be used in policies and CIB nodes to generate multiple outputs from a single entry.

**Policy Interface (PI):** The Policy Interface exposes a set of programming function calls that NEAT components may invoke to make requests to the NEAT Policy Manager. The JSON format was selected for the Policy Interface to achieve a decoupling of the PM from the rest of the NEAT System. The PM components described in this section may become optional or may even be executed outside the host using them in a future version of the NEAT System, e.g., for resource-constrained mobile devices.

The communication interface between the Policy Manager and other NEAT components is currently implemented using Unix domain sockets. The socket `neat_pm_socket` is used to receive requests from the application and serve back a list of connection candidates to the NEAT Logic. The application request is passed as a string containing a JSON object with a set of NEAT properties, as given in Listing 24. The request should at least contain a destination DNS name and port, or a destination IP, destination port, and a local interface and IP.

Similarly, the PM response contains a configurable number of candidates encoded as a JSON list, wherein each list element is a JSON object containing the properties associated with a particular candidate. The list is ordered by the total score of all candidate properties.

Two additional Unix sockets, `neat_pib_socket` and `neat_cib_socket`, are available for adding new policies and CIB nodes to the PIB and CIB respectively (in addition to the filesystem interface).

Additionally, the PM exposes a REST API intended to allow external applications, such as SDN controllers, to query the contents of the PIB/CIB and to populate these with new entries. If this optional API is started, the PM starts listening for HTTP connections on a predefined port (45888 by default). Applications may then access the addresses listed in Table 3 using HTTP's GET/PUT semantics (using JSON):

**Some examples of the operation:** An example of the current Policy Manager workflow is given below. Please refer to Fig. 11.

1. The Policy Manager receives a query from the NEAT Logic through the Policy Interface. The query contains a list of properties (compiled in JSON format) originating from the application requirements and the NEAT Logic (e.g., DNS lookup result) (see step 1. in Fig. 11). An example is provided in Listing 24.



2. The PM queries the PIB for profiles, which are used to map high level properties, such as `bulk_transfer|True`, to one or more concrete properties, such as `interface|eth0`, and `wireless|False` (step 2. in Fig. 11).
3. The PM performs a CIB lookup and receives an initial set of transport option candidates which – as far as possible – fulfil the query properties (i.e., by filtering the CIB by available interface types, path characteristics over relevant interfaces). Each candidate is a `PropertyArray` containing the matched query properties as well as the associated properties of the potential connection, e.g., supported interface features, TCP variants, cached remote endpoint capabilities, (step 3. in Fig. 11). See also § 3.4.3 for more details.
4. For each candidate the PM performs a PIB lookup. The lookup yields a set of candidates to which any matching policies have been applied (step 4. in Fig. 11). Each policy may extend the candidate `PropertyArray` with additional (optional or mandatory) properties defined by the policy. Any candidate which contradicts any mandatory policy properties is eliminated by the PM. See § 3.4.2 for more details.
5. The PM ranks the candidates based on the total score of the individual candidate properties. Finally the ranked candidate list is returned to the NEAT Logic (step 5. in Fig. 11). This information is used by the NEAT Logic to carry out the final connection selection.

**Provided Transport Service Feature(s):** While it does not actively offer any Transport Service Feature, this building block is indirectly involved in providing many Transport Service Features by other building blocks—e.g., Selection of a secure interface, NEAT flow delay budget, NEAT flow low latency, etc. can be mapped to policy attributes and combined into rules to define appropriate policies.

**Related building blocks:**

- PIB (§ 3.4.2).
- CIB (§ 3.4.3).

### 3.4.2 Policy Information Base (PIB)

This building block defines the PIB repository that stores policies in the NEAT System and is accessed by the Policy Manager.

**NEAT Policies:** Recall that a candidate is an object containing an arbitrary number of NEAT properties. NEAT policies are objects which append/update NEAT properties to connection *candidates* in order to enforce or express the desire for a specific behaviour. Each policy contains two essential keys: `match` and `properties`, as well as a number of meta attributes such as `uid`, `priority` or `description`. Policies are stored in the PIB repository.

Essentially NEAT policies are based on the following logic:

```
MATCH <set of properties to match against> → PROPERTIES <set of mandatory or optional properties appended to flow candidate>.
```

The list of match conditions always implies that conditions are evaluated by performing an AND operation, i.e., all conditions must be true for the policy to be triggered. To define an OR relationship between conditions, multiple policies with a different set of conditions must be created.

Table 4: Policy object keys.

Key	Description
<code>uid</code>	unique identifier for the policy (default: random UID)
<code>priority</code>	policies with the lowest priority are applied first (default: 0.0)
<code>description</code>	human readable description (default: "")
<code>policy_type</code>	either <code>policy</code> (default) or <code>profile</code>
<code>replace_matched</code>	if <code>true</code> the matched properties will be removed from the candidate prior to applying the policy properties (default: <code>false</code> )
<code>match</code>	properties that trigger the policy (default: <code>any</code> )
<code>properties</code>	policy properties which will update or be appended to the candidate (default: null set)

Table 4 provides a description of the policy object keys. If any of the above keys is omitted, its default value will be used.

- The value of the `match` key is an object which contains a set of properties that trigger the policy. A policy is triggered only if *all* of the match properties are contained within the input candidate, i.e., the set of candidate properties is a superset of the set of match properties. If the match key is empty or missing, the policy will match *any* candidate.
- The value of the `properties` key is a `PropertyArray` object (§ 3.4.1) which contains a set of properties that should be appended to the connection candidate, *if* the policy is triggered. To this end, the candidate properties are updated/appended by the *policy properties* specified in the properties object. Specifically, a candidate property will be *updated* (see *NEAT operators* in § 3.4.1) by a policy property if it has the same key. If a policy property's key is not already in the set of candidate properties the property will be appended to the set.

The NEAT Policy Manager uses JSON to encode policies. An example of a simple JSON policy is given in Listing 25.

```
1 {
2   "description":"TCP specific low latency profile",
3   "uid":"low_latency_tcp",
4   "priority": 3,
5   "replace_matched": true,
6   "match": {
7     "low_latency": {"value": true},
8     "transport": {"value": "TCP"}
9   },
10  "properties": {"transport": { "value": "TCP", "precedence": 2},
11                "SO/SOL_SOCKET/TCP_NODELAY": { "value": 1, "precedence": 1}}
12 }
```

Listing 25: JSON encoded policy.

**NEAT Profiles:** Profiles exhibit an identical behaviour to policies but are applied earlier — at the beginning of the workflow, before the CIB lookup (see Figure 11). The motivation is to allow the system to expand/map abstract properties (e.g., low latency, reliable transport), from a API request to one

or more concrete properties appropriate to the host on which the PM is running. Logically profiles belong to the PIB.

**PIB Scope:** We consider the following policy categories, each forming a separate *PIB domain*:

- *Global NEAT*: generic set coming from NEAT operations. It can only change (extended in order not to break compatibility) at the next update of the NEAT System.
- *Application specific*: set by each application (e.g., Mozilla Firefox). It can only change at the next update or installation of the application which it belongs to.
- *Operating system specific*: set by OS (e.g., Linux, FreeBSD, ...). It can only change at the next update or installation of the OS.
- *Vendor specific*: set by the end-host manufacturer (e.g., a handset maker). It can only change at the next update of the firmware on the end host.

In the current PIB implementation policies are stored as files in a predefined directory. The default locations for the policy files on Linux/Unix are:

- /etc/neat/policy/OS/
- /etc/neat/policy/vendor/
- /etc/neat/policy/application/
- \$HOME/.neat/policy/application/

The path configuration may be changed by the administrator at startup of the PM.

The PM imports all files into the PIB and maintains data structure for lookups. Note that PIB files are not intended to be accessed directly by NEAT components, but only through the Policy Interface. Read/write access to the PIB/CIB information is implemented using standard OS user/access mechanisms, Unix sockets, or through the PM REST API.

Optionally, the internal state of the PM may be exposed through an entry in `procfs` when available in Linux-like Oses (or a similar interface in other Oses).

**Provided Transport Service Feature(s):** There are no specific Transport Service Features associated to this building block.

**Some examples of the operation:** A sample of the proposed policy file format is given in Listing 26.

```
1 {
2   "uid":"low_latency",
3   "description":"low latency profile",
4   "priority": 1,
5   "replace_matched": true,
6   "match":{
7     "low_latency": {
8       "precedence": 1,
9       "value": true
10  }
```

```
11  },
12  "properties":{
13    "RTT": {
14      "precedence": 1,
15      "value": {"start":0, "end":50},
16      "score": 5
17    },
18    "low_latency_interface": {
19      "value": true, "precedence": 1},
20    "is_wired_interface": {
21      "value": true, "precedence": 2}
22  }
23 }
```

Listing 26: Policy file example.

**Related building blocks:**

- NEAT Policy Manager (§ 3.4.1).

**3.4.3 Characteristics Information Base (CIB)**

The Characteristics Information Base (CIB) is a repository that stores information about hosts (e.g., available interfaces, supported protocols), connections (e.g., parameters used by previously established transport sessions, hosts currently communicating) and the network (e.g., path properties). CIB entries provide measured information, protocol details and capabilities about network entities used in the NEAT System, specifically in the policy decision phase.

The CIB is populated from multiple inputs generated by so-called *CIB sources*, such as components of the NEAT System or external applications. A *CIB source* is defined as any module which provides an input for the CIB in the correct format accepted by the CIB. Each CIB source generates *CIB nodes*, which are JSON data objects used to populate the CIB, comprised of a set of attributes for a given resource. CIB sources may populate the CIB through two methods: by creating files stored in a pre-defined folder in the host's filesystem, the location of which is OS dependent, or through a REST API listening on a port of the host.

Some mechanisms to populate the CIB are already implemented in OSes as statistics/measurement tools and will be made available as default CIB sources. Another class of CIB sources is provided by NEAT building blocks such as Happy Eyeballs (§ 3.3.1) which store discovered transport protocols and parameters supported along paths in the CIB for future reuse. External CIB sources may be provided as helper applications by device and OS vendors or third parties developing modules for active or passive measurements, statistics and metadata collection.

The PM uses a pull mechanism to access any information stored in the CIB whenever a new NEAT Flow is initiated.

In the initial CIB implementation, trust is managed by relying on existing OS roles and permissions: CIB sources are allowed to create and update files in the CIB repository folder as long as the OS user executing the task is allowed to write in that folder. Future versions of the CIB may switch to a more complex authentication and trust management method.

**CIB architecture:** Each entry or *row* of the CIB is comprised of an arbitrary number of *NEAT Properties*. Rows are composed from one or more *CIB nodes* represented as JSON objects — and, e.g., stored as files in the CIB directory or received through the REST API. These JSON objects are generated by various helper applications called *CIB sources*. A simple example of the CIB node defining a local interface is given in Listing 27.

```
1 {
2   "uid": "eth0",
3   "root": true,
4   "expire": -1,
5   "priority": 4,
6   "properties": {
7     "interface": {"value": "eth0", "precedence":2},
8     "capacity": {"value": 10000, "precedence":2},
9     "local_ip": {"value": "10.10.2.1", "precedence":2},
10    "is_wired": {"value": true, "precedence":2},
11    "MTU": {"value": {"start":50, "end":9000}}
12  }
13 }
```

Listing 27: Example CIB node, showing the CIB format.

Each CIB node includes a `uid` key and a set of `properties` to be included in the CIB row.

**Extending existing CIB Nodes:** Additionally, CIB sources may generate CIB nodes which reference and extend pre-existing CIB nodes. This feature is used to inject network characteristics collected by external sources (e.g., controllers). For example, the following two CIB nodes, in Listings 28 and 29 reference the node defined in Listing 27:

```
1 {
2   "uid": "eth0_remote_1",
3   "description": "information about remote endpoint 1",
4   "priority": 2,
5   "link": true,
6   "match" : [
7     {"uid": {"value": "eth0"}}
8   ],
9   "properties": {
10    "remote_ip": {"value": "8.8.8.8", "precedence":2, "score": 2},
11    "remote_port": {"value": "80", "precedence":2, "score": 1}
12  }
13 }
```

Listing 28: Example of CIB node referencing the node in Listing 27. Referencing is done via the `uid` value in the `match` attribute (`eth0`).

```
1 {
2   "uid": "eth0_remote_2",
3   "description": "information about remote endpoint 2",
4   "priority": 2,
5   "link": true,
```

```
6  "match" : [  
7    {"interface": {"value": "eth0"}, "local_ip": {"value": "10.10.2.1"}}  
8  ],  
9  "properties": {  
10     "remote_ip": {"value": "8.8.4.4.", "precedence":2, "score": 2},  
11  }  
12 }
```

Listing 29: Example of CIB node referencing the node in Listing 27. Referencing is done via the `interface` and `local_ip` properties in the `match` attribute (`eth0`).

If a CIB node contains a `link` attribute that is set to `true` the CIB will attempt to match any property or `uid` listed in the `match` attribute, and will append the new properties to the corresponding CIB node. Both examples above will be linked to the interface `eth0`, defined in the first CIB node. To reference multiple CIBs, the `match` attribute list can contain multiple JSON objects. The `priority` attribute is used to resolve overlapping properties.

Essentially, the CIB internally constructs a directed graph using all available CIB nodes. To generate the CIB rows, the graph is resolved starting at each node which has the attribute `root` set to `true` (see CIB node in Listing 27), generating paths (i.e., rows) by traversing the graph in the reverse direction of the edges.

Hence, from the above CIB nodes the CIB will generate the following two rows:

```
1  
2 1: {"interface": {"value": "eth0", "precedence":2}, "capacity": {"value": 10000, "  
   precedence":2}, "local_ip": {"value": "10.10.2.1", "precedence":2}, "is_wired":  
   {"value": true, "precedence":2}, "MTU": {"value": {"start":50, "end":9000}, "  
   remote_ip": {"value": "8.8.8.8", "precedence":2, "score": 2}, "remote_port": {"  
   value": "80", "precedence":2, "score": 1}}  
3  
4 2: {"interface": {"value": "eth0", "precedence":2}, "capacity": {"value": 10000, "  
   precedence":2}, "local_ip": {"value": "10.10.2.1", "precedence":2}, "is_wired":  
   {"value": true, "precedence":2}, "MTU": {"value": {"start":50, "end":9000}, "  
   remote_ip": {"value": "8.8.4.4.", "precedence":2, "score": 2}}
```

Listing 30: CIB rows generated from the CIB nodes defined in Listings 27 to 29.

**Extending CIB Rows:** Finally, CIB sources have the option of generating CIB nodes which augment existing CIB *rows* with additional properties. This can be useful to annotate specific CIB rows with historical information, e.g., collected from previous connections. To achieve this the `link` attribute is set to `false`:

```
1 {  
2  "uid": "historic info",  
3  "description": "appended properties to CIB rows matching all match properties.",  
4  "priority": 10,  
5  "timestamp": 1476104788,  
6  "link": false,  
7  "match": [  
8    { "interface": {"value": "eth0"},  
9      "local_ip": {"value": "10.10.2.1"},
```

```
10     "remote_ip": {"value": "8.8.8.8"}}
11   ],
12   "properties": [{
13     "remote_port": {"value": 8080, "precedence":1},
14     "local_port": {"value": 56674, "precedence":1},
15     "transport": {"value": "TCP", "precedence":1},
16     "cached": {"value": true, "precedence":2, "score":5},
17     "cache_ttl": {"value": 300, "precedence":1},
18     "cache_status": {"value": "connection_success", "precedence":2, "description":
19       "could be failed, NA, etc."}
20   ] }
```

Listing 31: Extending CIB rows: any row matching all properties in the `match` attribute will be augmented with the properties in the `properties` attribute. In this example row 1 in Listing 30.

For this example the CIB will match row 1 in Listing 30 (i.e., the one with the `remote_ip` property set to 8.8.8.8), and insert a new row which includes the additional properties.

```
1 3: {"interface": {"value": "eth0", "precedence":2}, "capacity": {"value": 10000, "precedence":2}, "local_ip": {"value": "10.10.2.1", "precedence":2}, "is_wired": {"value": true, "precedence":2}, "MTU": {"value": {"start":50, "end":9000}, "remote_ip": {"value": "8.8.8.8", "precedence":2, "score": 2}, "remote_port": {"value": 8080, "precedence":1}, "local_port": {"value": 56674, "precedence":1}, "transport": {"value": "TCP", "precedence":1}, "cache_ttl": {"value": 300, "precedence":1}, "cached_connection_status": {"value": "success", "precedence":2, "score":5}}
```

Listing 32: CIB row 1 extended by the CIB node defined in Listing 31.

CIB nodes are *dynamic*, in the sense that they can be continuously updated—e.g., by a module monitoring the local OS capabilities or an external SDN controller. Further, whenever a new NEAT Flow is established, the Happy Eyeballs module caches the properties of the selected candidate by generating a corresponding CIB node referencing a local interface. These cache entries expire after a fixed time once the connection has been closed (time to live, TTL). The value of the TTL for remote CIB entries will be defined based on the initial experiences with the PM.

**CIB lookup workflow:** CIB lookup requests are arrays of NEAT properties including the address of the NEAT destination. Within the PM the request will contain the application requirements as well as any properties appended by the profile PIB. For each lookup request, the CIB compares the properties in the request against the properties contained in each row of the CIB. The lookup is successful whenever *all* required (precedence: 2) request properties match the properties in a CIB row. In this case, the properties of the corresponding CIB row are appended to the request, and the resulting property array becomes a candidate.

As a result, after the CIB lookup, the PM receives a list of candidates which fulfill the application requirements as well as the constraints and attributes imposed by the local system and the attached networks stored in the CIB. Candidates are ranked by the sum of the scores of the individual properties. However, only candidate properties which have been evaluated (i.e., matched) during the lookup are used to calculate the score, in order to ensure that properties about which nothing is known to the CIB

do not influence the ranking. Figure 4 illustrates how the CIB lookup process fits within the overall PM workflow.

Examples of NEAT-provided CIB sources include:

- Statistics and metadata provided by the operating system (e.g., network interface, socket statistics, battery drain, etc.).
- Statistics about path support from completed transport sessions and Happy Eyeballs (e.g., transport support and IP version).
- Path characteristics derived from various passive and active measurement techniques (e.g., network controller, network probes).
- Interface metadata (e.g., signal strength, type).

While some of these CIB sources have already been implemented (interface stats, network controller path information), additional CIB sources may be developed as NEAT is deployed in new scenarios.

**Provided Transport Service Feature(s):** There are no specific Transport Service Features associated to this building block.

**Related building blocks:**

- NEAT Policy Manager (§ 3.4.1).

## 4 NEAT reference material

A key element to a successful implementation project is proper documentation, ease of use of the code, and reference material like tutorials and pre-packaged examples. The NEAT consortium believes that the public community—e.g., network application developers, middleware developers, and students—should be able to learn how to use the NEAT library in the most simple and straightforward way. Therefore, the NEAT code is complemented by online documentation as well as compiling/installation files and instructions; also, API reference documentation as well as a tutorial have been made publicly available. In addition, a series of code examples can be found throughout this document with regards to each transport component. The NEAT documentation is maintained up-to-date as the library evolves using the *readthedocs* Content Management System [8].

### 4.1 NEAT tutorial

A tutorial on how to use the NEAT library, explaining the main concepts used in NEAT (e.g., flows, contexts, properties, etc.) and including a step-by-step minimal client and server example is available at <http://neat.readthedocs.io/en/latest/tutorial.html>. Section 2.1 of this document provides a copy of such tutorial at the time of writing.



Table 5: OS distributions used for NEAT VMs.

Operating System	Version
Ubuntu	16.04 LTS
Debian	9.0 (stretch)
FreeBSD	11

## 4.2 Additional online documentation

The NEAT API reference documentation describes a set of API functions exposed to the user. It also explains the callback mechanism and the related choices used in NEAT, as well as the error codes, the optional arguments, and properties. It is available at <http://neat.readthedocs.io/en/latest/index.html> and a to-date instance of it is included in Appendix B of this document. A reference to NEAT's coding style can be found at <http://neat.readthedocs.io/en/latest/internal/codingstyle.html>. It is based on the coding style used in the Linux kernel.

## 4.3 Virtual machines

To enable interested developers to quickly evaluate NEAT, the NEAT library, including installation instructions, sample test applications and documentation, is provided within ready-to-deploy virtual machine (VM) environments. These NEAT VMs support both VMware and VirtualBox hypervisors and provide a self-contained environment showcasing the NEAT library. Each VM contains a pre-configured version of the NEAT library, built from a recent snapshot of the Github repository. NEAT VMs are publicly available to download via <https://www.neat-project.org/resources/> to facilitate usage and adoption by interested developers. Users can evaluate and test NEAT in a safe and supported machine and they can clone it as many times as needed to simulate multiple hosts in complex topologies. The list of OS distributions (Long Term Support versions are preferred) available through the NEAT VM repository at the time of writing is listed in Table 5.

Each VM contains a recent binary version of the NEAT library that has successfully passed the automated Buildbot tests<sup>14</sup>. The Policy Manager is configured to run as a daemon, using a minimal set of pre-configured policies and CIB entries. In addition, the VMs contain the NEAT sources in the directory `~/neat` linked to the NEAT Github repository enabling users to available to pull updated development code or push contributions. The images also include all package dependencies required to build and debug NEAT, as well as compiler suites and example NEAT code.

The process of configuring and updating the VM installation of NEAT on Debian-based systems has been automated using Ansible<sup>15</sup>.

## 5 Conclusions

This deliverable presents the final version of a prototype implementation of the NEAT Core Transport System. The work realises NEAT as a user-space library. The components of the NEAT Transport Services have already been identified in Deliverable D1.3 [41], they provide an alternative to the legacy sockets API. The NEAT library allows applications to be developed that request the service they need

<sup>14</sup>These are described in detail in Deliverable D4.2 [14].

<sup>15</sup><https://github.com/NEAT-project/neat/tree/master/debian/ansible>

using a new approach that is agnostic to the specific choice of transport protocol and protocol parameters. This not only allows applications to be developed that can take advantage of common protocol machinery, but it also eases introduction of new network mechanisms and transport protocols.

The document reflects the final status of the implementation at the end of WP2 (milestone MS12), describing each of the components in the NEAT library, together with operation of the system as a whole and a set of examples of usage. It reflects the design and development efforts undertaken in Tasks 2.1 and 2.2. It also includes the final adjustment of the code in Task 2.3 to ensure the realised services and APIs are consistent with those defined in Task 1.4 of WP1. Further work to maintain the prototype and support analysis of performance and functionalities may continue within the NEAT project as a part of WP4, but are beyond the scope of this document.

Section 1 provides a short overview of the NEAT architecture (§ 1.1) and identified the Transport Services provided by the NEAT User API (§ 1.2). It introduced the core transport components required to provide such Transport Services (§ 1.3) and provided a step-by-step summary of using the NEAT workflow. This presented operation from the moment a connection setup attempt is made by an application using the API, to the point where a connection handle is returned to the application (§ 1.4), allowing the application to use a callback interface to send and receive network data. It also includes an overview of connection teardown when the application has no further use for the network connection (§ 1.5).

Section 2 presents a detailed tutorial on the new approach to developing applications programs using the NEAT User API. It contains a full example of a client/server application using NEAT (§ 2.1). In describing use of the API it summarised how this eases network programming in contrast to the traditional Berkeley socket API (§ 2.2), allowing a more streamlined code in network applications: Built-in functionalities underneath the NEAT User API that are transparent to the application programmer and supported across multiple OS platforms (e.g., name resolution, transport protocol selection, and fallback mechanisms); and the call-back approach simplifies network programming. A set of bindings for the Python scripting language were also developed and these are discussed in § 2.3.

Section 3 provides a detailed review of each of the core transport functions. This is structured into four classes of building blocks: Framework, Transport, Selection and Policy according to the architecture documented in Deliverable D1.1 [20]. The operation of each component is explained, as implemented in the final version of the NEAT library, and their internal dependencies are identified. Examples of use or operation has been provided together with code snippets in the C language.

Finally, Section 4 summarises the documentation developed for the NEAT core transport system. This includes an online NEAT tutorial, the NEAT API reference, and VMs pre-packaged with NEAT support. In summary, thanks to their platform- and protocol-independent nature, the NEAT core transport system and its User API as outlined in this document provide a simple, flexible, and easy way for application programmers to take advantage of the transport services provided by Internet transport protocols.

## References

- [1] Filestreamer. [Online]. Available: <https://github.com/oystedal/filestreamer>
- [2] libre — a portable and generic library for real-time communication. [Online]. Available: <http://creytiv.com/re.html>
- [3] libuv library. [Online]. Available: <http://libuv.org/>
- [4] Node.js. [Online]. Available: <https://nodejs.org/>
- [5] Openssl bio\_s\_mem manual. [Online]. Available: [https://wiki.openssl.org/index.php/Manual: BIO\\_s\\_mem\(3\)](https://wiki.openssl.org/index.php/Manual: BIO_s_mem(3))
- [6] OpenSSL library. [Online]. Available: <https://www.openssl.org/>
- [7] RAWRTC — webrtc and ortc library. [Online]. Available: <http://creytiv.com/re.html>
- [8] Read the docs. [Online]. Available: <https://readthedocs.org>
- [9] socket.io. [Online]. Available: <https://socket.io>
- [10] TCP increased security working group. [Online]. Available: <https://datatracker.ietf.org/wg/tcpinc/charter/>
- [11] usrsctp — a portable sctp userland stack. [Online]. Available: <https://github.com/sctplab/usrsctp/>
- [12] D. Balenson, “Privacy Enhancement for Internet Electronic Mail: Part III: Algorithms, Modes, and Identifiers,” RFC 1423 (Historic), Internet Engineering Task Force, Feb. 1993. [Online]. Available: <http://www.ietf.org/rfc/rfc1423.txt>
- [13] Z. Bozakov, S. Mangiante, C. Benet, A. Brunstrom, R. Santos, A. Kassler, and D. Buckley, “A NEAT framework for enhanced end-host integration in SDN environments,” in *IEEE Conference on Network Function Virtualization and Software Defined Networks (IEEE NFV-SDN)*, Berlin, Nov. 2017, accepted for publication, to appear.
- [14] Z. Bozakov, A. Brunstrom, D. Damjanovic, G. Fairhurst, A. E. Hansen, T. Jones, N. Khademi, A. Petlund, D. Ros, T. Rozensztrauch, M. I. Sánchez Bueno, D. Stenberg, M. Tüxen, and F. Weinrank, “Final Version of NEAT-based Tools,” The NEAT Project (H2020-ICT-05-2014), Deliverable D4.2, Sep. 2017.
- [15] V. Cerf, Y. Dalal, and C. Sunshine, “Specification of Internet Transmission Control Program,” RFC 675 (Historic), Internet Engineering Task Force, Dec. 1974, obsoleted by RFC 7805. [Online]. Available: <http://www.ietf.org/rfc/rfc675.txt>
- [16] Y. Cheng, J. Chu, S. Radhakrishnan, and A. Jain, “TCP Fast Open,” RFC 7413 (Experimental), Internet Engineering Task Force, Dec. 2014. [Online]. Available: <http://www.ietf.org/rfc/rfc7413.txt>
- [17] T. Dierks and E. Rescorla, “The Transport Layer Security (TLS) Protocol Version 1.2,” RFC 5246 (Proposed Standard), Internet Engineering Task Force, Aug. 2008, updated by RFCs 5746, 5878, 6176, 7465, 7507, 7568, 7627, 7685, 7905, 7919. [Online]. Available: <http://www.ietf.org/rfc/rfc5246.txt>

- [18] ECMA. (2013) The JSON data interchange format. [Online]. Available: <http://www.ecma-international.org/publications/files/ECMA-ST/ECMA-404.pdf>
- [19] G. Fairhurst, B. Trammell, and M. Kuehlewind, “Services Provided by IETF Transport Protocols and Congestion Control Mechanisms,” RFC 8095 (Informational), Internet Engineering Task Force, Mar. 2017. [Online]. Available: <http://www.ietf.org/rfc/rfc8095.txt>
- [20] G. Fairhurst, T. Jones, Z. Bozakov, A. Brunstrom, D. Damjanovic, T. Eckert, K. R. Evensen, K.-J. Grinnemo, A. F. Hansen, N. Khademi, S. Mangiante, P. McManus, G. Papastergiou, D. Ros, M. Tüxen, E. Vyncke, and M. Welzl, “NEAT Architecture,” The NEAT Project (H2020-ICT-05-2014), Deliverable D1.1, Dec. 2015. [Online]. Available: <https://www.neat-project.org/publications/>
- [21] K.-J. Grinnemo, A. Brunstrom, P. Hurtig, N. Khademi, and Z. Bozakov, “Happy Eyeballs for Transport Selection,” Internet Draft draft-grinnemo-taps-he, Jun. 2017, work in progress. [Online]. Available: <https://tools.ietf.org/html/draft-grinnemo-taps-he>
- [22] K.-J. Grinnemo, Z. Bozakov, A. Brunstrom, D. Damjanovic, K. Evensen, G. Fairhurst, A. Hansen, D. Hayes, P. Hurtig, N. Khademi, S. Mangiante, A. Mohideen, M. Rajiullah, D. Ros, I. Rüngeler, M. I. Sánchez Bueno, R. Santos, R. Secchi, T. Tangenes, M. Tüxen, F. Weinrank, and M. Welzl, “Initial Report on the Extended Transport System,” The NEAT Project (H2020-ICT-05-2014), Deliverable D3.1, Aug. 2016.
- [23] K.-J. Grinnemo, A. Brunstrom, G. Fairhurst, D. Hayes, P. Hurtig, N. Khademi, D. Ros, I. Rüngeler, M. Tüxen, F. Weinrank, and M. Welzl, “Final Report on Transport Protocol Enhancements,” The NEAT Project (H2020-ICT-05-2014), Deliverable D3.2, Mar. 2017.
- [24] M. Handley, V. Jacobson, and C. Perkins, “SDP: Session Description Protocol,” RFC 4566 (Proposed Standard), Internet Engineering Task Force, Jul. 2006. [Online]. Available: <http://www.ietf.org/rfc/rfc4566.txt>
- [25] D. Hayes, D. Ros, A. Petlund, and I. Ahmed, “A framework for less than best effort congestion control with soft deadlines,” in *Proceedings of IFIP Networking*, Stockholm, Jun. 2017. [Online]. Available: <http://dl.ifip.org/db/conf/networking/networking2017/1570334752.pdf>
- [26] S. Islam, M. Welzl, K. Hiorth, D. Hayes, Ø. Dale, G. Armitage, and S. Gjessing, “Single-path TCP Congestion Control Coupling,” Technical Report 459, Department of Informatics, University of Oslo, Jan. 2017.
- [27] B. Kaliski, “Privacy Enhancement for Internet Electronic Mail: Part IV: Key Certification and Related Services,” RFC 1424 (Historic), Internet Engineering Task Force, Feb. 1993. [Online]. Available: <http://www.ietf.org/rfc/rfc1424.txt>
- [28] S. Kent, “Privacy Enhancement for Internet Electronic Mail: Part II: Certificate-Based Key Management,” RFC 1422 (Historic), Internet Engineering Task Force, Feb. 1993. [Online]. Available: <http://www.ietf.org/rfc/rfc1422.txt>
- [29] N. Khademi, D. Ros, M. Welzl, Z. Bozakov, A. Brunstrom, G. Fairhurst, K.-J. Grinnemo, D. Hayes, P. Hurtig, T. Jones, S. Mangiante, M. Tüxen, and F. Weinrank, “NEAT: A Platform- and Protocol-Independent Internet Transport API,” *IEEE Communications Magazine*, vol. 55, no. 6, pp. 46–54, Jun. 2017.

- [30] N. Khademi, Z. Bozakov, A. Brunstrom, D. Damjanovic, K. R. Evensen, G. Fairhurst, K.-J. Grinnemo, T. Jones, S. Mangiante, A. Petlund, D. Ros, D. Stenberg, M. Tüxen, F. Weinrank, and M. Welzl, “Core Transport System, with both Low-level and High-level Components,” NEAT Project (H2020-ICT-05-2014), Deliverable D2.2, Mar. 2017. [Online]. Available: <https://www.neat-project.org/publications/>
- [31] J. Linn, “Privacy Enhancement for Internet Electronic Mail: Part I: Message Encryption and Authentication Procedures,” RFC 1421 (Historic), Internet Engineering Task Force, Feb. 1993. [Online]. Available: <http://www.ietf.org/rfc/rfc1421.txt>
- [32] R. Mahy, P. Matthews, and J. Rosenberg, “Traversal Using Relays around NAT (TURN): Relay Extensions to Session Traversal Utilities for NAT (STUN),” RFC 5766 (Proposed Standard), Internet Engineering Task Force, Apr. 2010. [Online]. Available: <http://www.ietf.org/rfc/rfc5766.txt>
- [33] G. Papastergiou, K.-J. Grinnemo, A. Brunstrom, D. Ros, M. Tüxen, N. Khademi, and P. Hurtig, “On the cost of using Happy Eyeballs for transport protocol selection,” in *Proceedings of the 2016 Applied Networking Research Workshop (ANRW)*, Berlin, Jul. 2016, pp. 45–51.
- [34] E. Rescorla and N. Modadugu, “Datagram Transport Layer Security Version 1.2,” RFC 6347 (Proposed Standard), Internet Engineering Task Force, Jan. 2012, updated by RFCs 7507, 7905. [Online]. Available: <http://www.ietf.org/rfc/rfc6347.txt>
- [35] J. Rosenberg, “Interactive Connectivity Establishment (ICE): A Protocol for Network Address Translator (NAT) Traversal for Offer/Answer Protocols,” RFC 5245 (Proposed Standard), Internet Engineering Task Force, Apr. 2010, updated by RFC 6336. [Online]. Available: <http://www.ietf.org/rfc/rfc5245.txt>
- [36] J. Rosenberg, R. Mahy, P. Matthews, and D. Wing, “Session Traversal Utilities for NAT (STUN),” RFC 5389 (Proposed Standard), Internet Engineering Task Force, Oct. 2008, updated by RFC 7350. [Online]. Available: <http://www.ietf.org/rfc/rfc5389.txt>
- [37] R. Stewart, M. Tuexen, S. Loreto, and R. Seggelmann, “Stream schedulers and user message interleaving for the stream control transmission protocol,” Working Draft, IETF Secretariat, Internet-Draft draft-ietf-tsvwg-sctp-ndata-12, Aug. 2017. [Online]. Available: <http://www.ietf.org/internet-drafts/draft-ietf-tsvwg-sctp-ndata-12.txt>
- [38] F. Weinrank, K.-J. Grinnemo, Z. Bozakov, A. Brunstrom, T. Dreibholz, G. Fairhurst, P. Hurtig, N. Khademi, and M. Tüxen, “A NEAT Way to Browse the Web,” in *Proceedings of the 2017 Applied Networking Research Workshop (ANRW)*, Prague, Jul. 2017. [Online]. Available: <https://irtf.org/anrw/2017/anrw17-final13.pdf>
- [39] M. Welzl, M. Tuexen, and N. Khademi, “On the usage of transport service features provided by IETF transport protocols,” Internet Draft draft-ietf-taps-transports-usage, Aug. 2017, work in progress. [Online]. Available: <https://tools.ietf.org/html/draft-ietf-taps-transports-usage>
- [40] M. Welzl, A. Brunstrom, D. Damjanovic, K. Evensen, T. Eckert, G. Fairhurst, N. Khademi, S. Mangiante, A. Petlund, D. Ros, and M. Tüxen, “First Version of Services and APIs,” The NEAT Project (H2020-ICT-05-2014), Deliverable D1.2, Mar. 2016. [Online]. Available: <https://www.neat-project.org/publications/>

- [41] M. Welzl, D. Damjanovic, T. Jones, M. Tüxen, and F. Weinrank, “Final Version of Services and APIs,” The NEAT Project (H2020-ICT-05-2014), Deliverable D1.3, Sep. 2017.
- [42] D. Wing and A. Yourtchenko, “Happy Eyeballs: Success with Dual-Stack Hosts,” RFC 6555 (Proposed Standard), Internet Engineering Task Force, Apr. 2012. [Online]. Available: <http://www.ietf.org/rfc/rfc6555.txt>

## A NEAT Terminology

This appendix defines terminology used to describe NEAT. These terms are used throughout this document.

**Application** An entity (program or protocol module) that uses the transport layer for end-to-end delivery of data across the network (this may also be an upper layer protocol or tunnel encapsulation). In NEAT, the application data is communicated across the network using the NEAT User API either directly, or via middleware or a NEAT Application Support API on top of the NEAT User API.

**Characteristics Information Base (CIB)** The entity where path information and other collected data from the NEAT System is stored for access via the NEAT Policy Manager.

**NEAT API Framework** A callback-based API in NEAT. Once the NEAT base structure has started, using this framework an application can request a connection (create NEAT Flow), communicate over it (write data to the NEAT Flow and read received data from the NEAT Flow) and register callback functions that will be executed upon the occurrence of certain events.

**NEAT Application Support Module** Example code and/or libraries that provide a more abstract way for an application to use the NEAT User API. This could include methods to directly support a middleware library or an interface to emulate the traditional Socket API.

**NEAT Component** An implementation of a feature within the NEAT System. An example is a “Happy Eyeballs” component to provide Transport Service selection. Components are designed to be portable (e.g. platform-independent).

**NEAT Diagnostics and Statistics Interface** An interface to the NEAT System to access information about the operation and/or performance of system components, and to return endpoint statistics for NEAT Flows.

**NEAT Flow** A flow of protocol data units sent via the NEAT User API. For a connection-oriented flow, this consists of the PDUs related to a specific connection.

**NEAT Flow Endpoint** The NEAT Flow Endpoint is a NEAT structure that has a similar role to the Transmission Control Block (TCB) in the context of TCP. This is mainly used by the NEAT Logic to collect the information about a NEAT Flow.

**NEAT Framework** The Framework components include supporting code and data structures needed to implement the NEAT User Module. They call other components to perform the functions required to select and realise a Transport Service. The NEAT User API is an important component of the NEAT Framework; other components include diagnostics and measurement.

**NEAT Logic** The NEAT Logic is at the core of the NEAT System as part of the NEAT Framework components and is responsible for providing functionalities behind the NEAT User API.

**NEAT Policy Manager** Part of the NEAT User Module responsible for the policies used for service selection. The Policy Manager is accessed via the (user-space) Policy Interface, portable across platforms. An implementation of the NEAT Policy Manager may optionally also interface to kernel functions or implement new functions within the kernel (e.g. relating to information about a specific network interface or protocols).

**NEAT Selection** Selection components are responsible for choosing an appropriate transport endpoint and a set of transport components to create a Transport Service Instantiation. This utilises information passed through the NEAT User API, and combines this with inputs from the NEAT Policy Manager to identify candidate services and test the suitability of the candidates to make a final selection.

**NEAT Signalling and Handover** Signalling and Handover components enable optional interaction with remote endpoints and network devices to signal the service requested by a NEAT Flow, or to interpret signalling messages concerning network or endpoint capabilities for a Transport Service Instantiation.

**NEAT System** The NEAT System includes all user-space and kernel-space components needed to realise application communication across the network. This includes all of the NEAT User Module, and the NEAT Application Support Module.

**NEAT User API** The API to the NEAT User Module through which application data is exchanged. This offers Transport Services similar to those offered by the Socket API, but using an event-driven style of interaction. The NEAT User API provides the necessary information to allow the NEAT User Module to select an appropriate Transport Service. This is part of the NEAT Framework group of components.

**NEAT User Module** The set of all components necessary to realise a Transport Service provided by the NEAT System. The NEAT User Module is implemented in user space and is designed to be portable across platforms. It has five main groupings of components: Selection, Policy (i.e. the Policy Manager and its related information bases and default values), Transport, Signalling and Handover, and the NEAT Framework. The NEAT User Module is a subset of a NEAT System.

**Policy Information Base (PIB)** The rules used by the NEAT Policy Manager to guide the selection of the Transport Service Instantiation.

**Policy Interface (PI)** The interface to allow querying of the NEAT Policy Manager.

**Stream** A set of data blocks that logically belong together, such that uniform network treatment would be desirable for them. A stream is bound to a NEAT Flow. A NEAT Flow contains one or more streams.

**Transport Address** A transport address is defined by a network-layer address, a transport-layer protocol, and a transport-layer port number.

**Transport Feature** Short for Transport Service Feature.

**Transport Service** A set of end-to-end features provided to users, without an association to any given framing protocol, which provides a complete service to an application. The desire to use a specific feature is indicated through the NEAT User API.

**Transport Service Feature** A specific end-to-end feature that the transport layer provides to an application. Examples include confidentiality, reliable delivery, ordered delivery and message-versus-stream orientation.



**Transport Service Instantiation** An arrangement of one or more transport protocols with a selected set of features and configuration parameters that implements a single Transport Service. Examples include: a protocol stack to support TCP, UDP, or SCTP over UDP with the partial reliability option.

## B NEAT API Reference

### B.1 Optional arguments

Some of the functions in the NEAT API, such as `neat_open`, `neat_read` and `neat_write`, take optional arguments. These are sometimes used to pass optional arguments to functions, and sometimes used to return additional values from the function. Optional arguments are passed as an array of the struct `neat_tlv` and an integer specifying the length of this array. `neat_tlv` is defined as follows:

```
1 struct neat_tlv {
2     neat_tlv_tag tag;
3     neat_tlv_type type;
4
5     union {
6         int integer;
7         char *string;
8         float real;
9     } value;
10 };
```

Listing 33: `neat_tlv` struct.

An optional argument takes the form of a tag name, a type, and a value of either a string, integer or a floating point number. The tag specifies which optional argument the value belongs to, and the type asserts the type of the value passed as this argument. An error will be raised if the type is different than what the function expects.

You may either work with this struct directly, or you may use the preprocessor macros explained later in this document.

#### B.1.1 Specifying no optional arguments

To specify no optional arguments, simply pass `NULL` as the `optargs` parameter and `0` as the `opt_count` parameter of the function.

#### B.1.2 Optional argument macros

These are the optional argument macros used in NEAT:

- `NEAT_OPTARGS_DECLARE (max)`: Declare the necessary variables to use the rest of these macros. Allocates (on the stack) an array of length `max` and an integer for storing the number of optional arguments specified. `NEAT_OPTARGS_MAX` may be used as the default array size.
- `NEAT_OPTARGS_INIT ()`: Initializes the variables declared by `NEAT_OPTARGS_DECLARE`. May also be used to reset the (number of) optional arguments back to `0`.

- `NEAT_OPTARG_INT (tag, value)`: Specify the tag and the value of an optional argument that takes an integer.
- `NEAT_OPTARG_FLOAT (tag, value)`: Specify the tag and the value of an optional argument that takes a floating point number.
- `NEAT_OPTARG_STRING (tag, value)`: Specify the tag and the value of an optional argument that takes a string.
- `NEAT_OPTARGS`: Represents the array of optional arguments. Specify this macro as the `optarg` parameter.
- `NEAT_OPTARG_COUNT`: Stores the number of the optional arguments specified so far with `NEAT_OPTARG_INT`, `NEAT_OPTARG_FLOAT` or `NEAT_OPTARG_STRING`. This count is reset by `NEAT_OPTARGS_INIT ()`. Specify this macro as the `opt_count` argument.

### B.1.3 Optional argument tags

These are the optional argument tags used in NEAT:

- `NEAT_TAG_STREAM_ID (integer)`: Specifies the ID of the stream which the data should be written to, or which stream the data was read from.
- `NEAT_TAG_STREAM_COUNT (integer)`: Specifies the number of stream to create. Only used with protocols that support multiple streams.
- `NEAT_TAG_FLOW_GROUP (integer)`: Specifies the flow group this flow belongs to.
- `NEAT_TAG_PRIORITY (float)`: Specifies the priority of this flow relative to other flows in the flow group.
- `NEAT_TAG_CC_ALGORITHM (string)`: Specifies the name of the (TCP) congestion control algorithm that will be used by this flow. A system default will be used if the specified algorithm is not available.

Currently unused tags:

- `NEAT_TAG_LOCAL_NAME`
- `NEAT_TAG_SERVICE_NAME`
- `NEAT_TAG_CONTEXT`
- `NEAT_TAG_PARTIAL_RELIABILITY_METHOD`
- `NEAT_TAG_PARTIAL_RELIABILITY_VALUE`
- `NEAT_TAG_PARTIAL_MESSAGE_RECEIVED`
- `NEAT_TAG_PARTIAL_SEQNUM`
- `NEAT_TAG_UNORDERED`
- `NEAT_TAG_UNORDERED_SEQNUM`

- NEAT\_TAG\_DESTINATION\_IP\_ADDRESS

Example:

```
1  NEAT_OPTARGS_DECLARE (NEAT_OPTARGS_MAX);
2  NEAT_OPTARGS_INIT ();
3  NEAT_OPTARG_INT (NEAT_TAG_STREAM_COUNT, 5);
4  neat_open(ctx, flow, "127.0.0.1", 8000, NEAT_OPTARGS, NEAT_OPTARGS_COUNT);
```

## B.2 Properties

A property in NEAT may either express a requirement or it may express a desire from the application with regards to the service provided by the transport layer.

Properties are represented as JSON objects. A set of properties may be contained within one JSON object. Below is an example of a JSON object with a single property:

```
1 {
2   "property_name": {
3     value: "property_value",
4     precedence: 1
5   }
6 }
```

Note that all examples of properties will be specified inside a JSON object.

Properties have a name, a value, and a precedence. A string is always used for the name of a property. The value of a property may be either a boolean, a string, an integer, a floating point number, an array, or an interval. Each property expects only one specific type.

The properties are sent to the Policy Manager (if present), which will return a list containing a list of candidates, which is ordered by how good the candidate matches the request from the application. Each candidate specifies a given setting for each property. NEAT will use the properties specified in each candidate when trying to set up a new connection.

Some properties are set by NEAT based on parameters to function calls. Other properties must be set manually with the `neat_set_property` function.

### B.2.1 Application property reference

These are properties that may be set by the application.

**transport** Type: Array

Specifies an array of transport protocols that may be used. An application may specify either one protocol with precedence 2, or multiple protocols with precedence 1.

Note: May not be queried with `neat_get_property` before execution of the `on_connected` callback. When querying this property, the returned value is a string describing the actual transport in use.

Note: Applications should avoid specifying the protocol(s) to use directly, and instead rely on the Policy Manager to make a decision on what protocol(s) to use based on other properties. The transport property should only be used for systems without a Policy Manager, or if the choice of transport protocol is strictly mandated by the application protocol.

```
1 {
2   "transport": [
```

```
3     {
4         "value": "SCTP",
5         "precedence": 1
6     },
7     {
8         "value": "TCP",
9         "precedence": 1
10    }
11 ]
12 }
```

Listing 34: Example 1: multiple protocols

```
1 {
2     "transport": [
3         {
4             "value": "UDP",
5             "precedence": 2
6         }
7     ]
8 }
```

Listing 35: Example 2: one protocol

Available protocols:

- SCTP
- SCTP/UDP (SCTP tunneled over UDP)
- TCP
- UDP
- UDP-Lite

**Security** Type: Boolean

Specifies whether the connection should be encrypted or not. With precedence 2, NEAT will only report the connection as established if it was able to connect and the (D)TLS handshake succeeds. With precedence 1, NEAT may still attempt to establish an unencrypted connection.

### B.2.2 Inferred properties

These are properties that are inferred during connection setup and subsequently sent to the Policy Manager. Applications should not set these directly.

**interfaces** Type: Array

Specifies a list of available interfaces that may be used for connections. The Policy Manager may not use all interfaces in this list.

This property is inferred during the `neat_open` call. Do not set this property manually.

**domain\_name** Type: String

Specifies the name of the remote endpoint to connect to with the `neat_open` call.

This property is inferred from the name parameter of `neat_open` call. Do not set this property manually.

**port** Type: Integer

This property is inferred from the `neat_open` and `neat_accept` calls. Do not set this property manually.

### B.3 Callbacks

Callbacks are used in NEAT to signal events to the application. They are used to inform the application when a flow is readable, writable, or an error has occurred.

Most callbacks have the following syntax:

```
1 neat_error_code
2 on_event_name(neat_flow_operations *ops)
3 {
4     return NEAT_OK; // or some error code
5 }
```

Listing 36: NEAT callback syntax.

The struct `neat_flow_operations` is defined as follows:

```
1 struct neat_flow_operations
2 {
3     void *userData;
4
5     neat_error_code status;
6     int stream_id;
7     struct neat_ctx *ctx;
8     struct neat_flow *flow;
9
10    neat_flow_operations_fx on_connected;
11    neat_flow_operations_fx on_error;
12    neat_flow_operations_fx on_readable;
13    neat_flow_operations_fx on_writable;
14    neat_flow_operations_fx on_all_written;
15    neat_flow_operations_fx on_network_status_changed;
16    neat_flow_operations_fx on_aborted;
17    neat_flow_operations_fx on_timeout;
18    neat_flow_operations_fx on_close;
19    neat_cb_send_failure_t on_send_failure;
20    neat_cb_flow_slowdown_t on_slowdown;
21    neat_cb_flow_rate_hint_t on_rate_hint;
22 };
```

Listing 37: `neat_flow_operations` struct.

- `userData`: Applications may freely store a pointer in this field.
- `status`: Reports any errors associated with the flow.

- `stream_id`: For flows that use explicit multi-streaming. Specifies which stream the event is related to, if any.
- `ctx`: Pointer to the context the flow belongs to.
- `flow`: Pointer to the flow on which the event happened.

Callbacks are set by assigning the function pointer to the struct passed to the callback and then calling `neat_set_operations`. A NULL pointer may be used to indicate that the callback should no longer be called.

### B.3.1 Example callback flow

For most applications it will be sufficient to use the following callback flow as in Figure 6. See the tutorial in Section 2.1 for more details.

### B.3.2 Callback reference

Here is a list of callbacks used by NEAT.

`on_connected`: Called whenever an outgoing connection has been established with `neat_open`, or an incoming connection has been established with `neat_accept`.

`on_error`: Called whenever an error occurs when processing the flow. Errors are considered critical.

`on_readable`: Called whenever the flow can be read from without blocking. NEAT does not permit blocking reads.

`on_writable`: Called whenever the flow can be written to without blocking. NEAT does not permit blocking writes.

`on_all_written`: Called when all previous data sent with `neat_write` has been completely written. Does not signal that the flow is writable. Applications may use this callback to re-enable the `on_writable` callback.

`on_network_status_changed`: Inform application that something has happened in the network. This also includes flow endpoints going up, which will subsequently trigger `on_connected` if that callback is set. *Only available when using SCTP.*

`on_aborted`: Called when the remote end aborts the flow. Available for flows using TCP or SCTP.

`on_timeout`: Called if sent data is not acknowledged within the time specified with `neat_change_timeout`. *Currently only available for TCP on Linux.*

`on_close`: Called when the graceful connection shutdown has completed. Only available when using SCTP or TCP. Note that when using TCP, this callback is called when the `close()` system call is made, as TCP implementations currently does not provide any more accurate way of signalling this.

`on_send_failure`: Defined as:

```
1 void
2 on_send_failure(struct neat_flow_operations *flowops, int context, const unsigned
   char *unsent)
3 {
4 }
```

Called to inform the application that the returned message `unsent` could not be transmitted. The failure reason as reported by the transport protocol is returned in the standard status code, as an

abstracted NEAT error code. If the message was tagged with a context number, it is returned in context. Only available for SCTP. Flows using TCP may use timeouts instead.

`on_slowdown`: Not currently implemented. Defined as:

```
1 void
2 on_slowdown(struct neat_flow_operations *ops, int ecn, uint32_t rate)
3 {
4 }
```

Inform the application that the flow has experienced congestion and that the sending rate should be lowered. If `rate` is non-zero, it is an estimate of the new maximum sending rate. `ecn` is a boolean indicating whether this notification was triggered by an ECN mark.

`on_rate_hint`: Not currently implemented.

Defined as:

```
1 void
2 on_rate_hint(struct neat_flow_operations *ops, uint32_t new_rate)
3 {
4 }
```

Called to inform the application that it may increase its sending rate. If `new_rate` is non-zero, it is an estimate of the maximum sending rate.

## B.4 Error codes

These are the error codes used by NEAT library.

- `NEAT_OK`: Signals that no error has occurred. Equals to 0.
- `NEAT_ERROR_WOULD_BLOCK`: Signals that the operation could not be performed because it would block the process. NEAT does not permit blocking operations.
- `NEAT_ERROR_BAD_ARGUMENT`: Signals that one or more arguments given to the function was invalid or incorrect. This also includes optional arguments.
- `NEAT_ERROR_IO`: Signals that an internal I/O operation in NEAT has failed.
- `NEAT_ERROR_DNS`: Signals that there was an error performing DNS resolution.
- `NEAT_ERROR_INTERNAL`: Signals that there was an error internally in NEAT.
- `NEAT_ERROR_SECURITY`: Signals that there was an error setting up an encrypted flow.
- `NEAT_ERROR_UNABLE`: Signals that NEAT is not able to perform the requested operation.
- `NEAT_ERROR_MESSAGE_TOO_BIG`: Signals that the provided buffer space is not sufficient for the received message.
- `NEAT_ERROR_REMOTE`: Signals that there was an error on the remote endpoint.
- `NEAT_ERROR_OUT_OF_MEMORY`: Signals that NEAT is not able to allocate enough memory to complete the requested operation.

## B.5 API functions

### B.5.1 `neat_init_ctx`

- Syntax:

```
1 struct neat_ctx *neat_init_ctx();
```

- Parameters: None.
- Return values: Returns a pointer to a newly allocated and initialized NEAT context. Returns NULL if a context could not be allocated.
- Remarks: None.
- Examples: None.
- See also: `neat_free_ctx` and `neat_new_flow`

### B.5.2 `neat_free_ctx`

- Syntax:

```
1 void neat_free_ctx(struct neat_ctx *ctx);
```

- Parameters: `ctx`: Pointer to the NEAT context to free.
- Return values: None
- Remarks: If there are any flows still kept in this context, those will be freed and closed as part of this operation.
- Examples: None.
- See also: `neat_close` and `neat_init_ctx` and `neat_new_flow`

### B.5.3 `neat_new_flow`

- Summary: Allocate and initialize a new NEAT flow.
- Syntax:

```
1 neat_flow *neat_new_flow(struct neat_ctx *ctx);
```
- Parameters: `ctx`: Pointer to a NEAT context.
- Return values: Returns a pointer to a new flow. Returns NULL on error.
- Remarks: None.
- Examples: None.
- See also: `neat_open` and `neat_close`



#### B.5.4 neat\_set\_property

- Summary: Set the properties of a NEAT flow.
- Syntax:

```
1 neat_error_code neat_set_property(  
2     struct neat_ctx *ctx,  
3     struct neat_flow *flow,  
4     const char *properties);
```

- Parameters:
  - ctx: Pointer to a NEAT context.
  - flow: Pointer to a NEAT flow.
  - properties: Pointer to a JSON-encoded string containing the flow properties.
- Return values: Returns NEAT\_OK if the properties were set successfully. Returns NEAT\_ERROR\_BAD\_ARGUMENT if the JSON-encoded string is malformed.
- Remarks: Properties are applied when a flow connects.
- Examples: None.
- See also: None.

#### B.5.5 neat\_get\_property

- Summary: Query the properties of a flow. Returns value only, not precedence.
- Syntax:

```
1 neat_error_code neat_get_property(struct neat_ctx *ctx,  
2                                 struct neat_flow *flow,  
3                                 const char *name,  
4                                 void *ptr,  
5                                 size_t *size);
```

- Parameters:
  - ctx: Pointer to a NEAT context.
  - flow: Pointer to a NEAT flow.
  - name: Name of the property to query.
  - ptr: Pointer to a buffer where the property value may be stored.
  - size: Pointer to an integer containing the size of the buffer pointed to by ptr. Updated to contain the size of the property upon return.
- Return values: Returns NEAT\_OK if the property existed and there was sufficient buffer space available. The size parameter is updated to contain the actual size. Returns NEAT\_ERROR\_MESSAGE\_TOO\_BIG if there was not sufficient buffer space. The size parameter is updated to contain the required buffer size. Returns NEAT\_ERROR\_UNABLE if the property does not exist.

- **Remarks:** Applications may pass 0 as the size parameter to query the size of the property.
- **Examples:**

```
1     size_t bufsize = 0;
2     char buffer = NULL;
3
4     if (neat_get_property(ctx, flow, "transport", buffer, &bufsize) ==
        NEAT_ERROR_MESSAGE_TOO_BIG) {
5         buffer = malloc(bufsize);
6         if (buffer && neat_get_property(ctx, flow, "transport", buffer, &
            bufsize) == NEAT_OK) {
7             printf("Transport: %s\n", buffer);
8         }
9         if (buffer)
10            free(buffer);
11    } else {
12        printf("\tTransport: Error: Could not find property \"transport\"\n");
13    }
```

- See also: Properties and `neat_set_property`

### B.5.6 `neat_open`

- **Summary:** Open a neat flow and connect it to a given remote name and port.
- **Syntax:**

```
1 neat_error_code neat_open(struct neat_ctx *ctx,
2                          struct neat_flow *flow,
3                          const char *name,
4                          uint16_t port,
5                          struct neat_tlv optional[],
6                          unsigned int opt_count);
```

- **Parameters:**
  - `ctx`: Pointer to a NEAT context.
  - `flow`: Pointer to a NEAT flow.
  - `name`: The remote name to connect to.
  - `port`: The remote port to connect to.
  - `optional`: An array containing optional parameters.
  - `opt_count`: The length of the array containing optional parameters.
- **Optional parameters:**
  - `NEAT_TAG_STREAM_COUNT` (integer): The number of streams to open, for protocols that supports multistreaming. Note that NEAT may automatically make use of multi-streaming for multiple NEAT flows between the same endpoints when this parameter is not used.

- NEAT\_TAG\_FLOW\_GROUP (integer): The group ID that this flow belongs to. For use with coupled congestion control.
- NEAT\_TAG\_PRIORITY (float): The priority of this flow relative to the other flows. Must be between 0.1 and 1.0.
- NEAT\_TAG\_CC\_ALGORITHM (string): The congestion control algorithm to use for this flow.
- Return values: Returns NEAT\_OK if the flow opened successfully. Returns NEAT\_ERROR\_OUT\_OF\_MEMORY if the function was unable to allocate enough memory.
- Remarks: Callbacks can be specified with neat\_set\_operations. The on\_connected callback will be invoked if the connection established successfully. The on\_error callback will be invoked if NEAT is unable to connect to the remote endpoint.
- Examples:  

```
1 neat_open(ctx, flow, "bsd10.fh-muenster.de", 80, NULL, 0);
```
- See also: neat\_read and Optional arguments

### B.5.7 neat\_accept

- Summary: Listen to incoming connections on a given port on one or more protocols.
- Syntax:  

```
1 neat_error_code neat_accept(struct neat_ctx *ctx,  
2                             struct neat_flow *flow,  
3                             uint16_t port,  
4                             struct neat_tlv optional[],  
5                             unsigned int opt_count);
```
- Parameters:
  - ctx: Pointer to a NEAT context.
  - flow: Pointer to a NEAT flow.
  - port: The local port to listen for incoming connections on.
  - optional: An array containing optional parameters.
  - opt\_count: The length of the array containing optional parameters.
- Optional parameters:
  - NEAT\_TAG\_STREAM\_COUNT (integer): The number of streams to accept, for protocols that supports multistreaming.
- Return values: Returns NEAT\_OK if NEAT is listening for incoming connections on at least one protocol. Returns NEAT\_ERROR\_UNABLE if there is no appropriate protocol available for the flow properties that was specified. Returns NEAT\_ERROR\_BAD\_ARGUMENT if flow is pointing to a flow that is already opened or listening for incoming connections. Returns NEAT\_ERROR\_BAD\_ARGUMENT if NEAT\_TAG\_STREAM\_COUNT is less than 1. Returns NEAT\_ERROR\_OUT\_OF\_MEMORY if the function was unable to allocate enough memory.

- **Remarks:** Callbacks can be specified with `neat_set_operations`. The `on_connected` callback will be invoked if the connection established successfully. The `on_error` callback will be invoked if NEAT is unable to connect to the remote endpoint. Which protocols to listen to is determined by the flow properties.

- **Examples:**

```
1 neat_accept(ctx, flow, 8080, NULL, 0);
```

- **See also:** `neat_open` and **Optional arguments**

### B.5.8 `neat_read`

- **Summary:** Read data from a neat flow.

Should only be called from within the `on_readable` callback specified with `neat_set_operations`.

- **Syntax:**

```
1 neat_error_code neat_read(struct neat_ctx *ctx,  
2                          struct neat_flow *flow,  
3                          unsigned char *buffer,  
4                          uint32_t amount,  
5                          uint32_t *actual_amount,  
6                          struct neat_tlv optional[],  
7                          unsigned int opt_count);
```

- **Parameters:**

- `ctx`: Pointer to a NEAT context.
- `flow`: Pointer to a NEAT flow.
- `buffer`: Pointer to a buffer where read data may be stored.
- `amount`: The size of the buffer pointed to by `buffer`.
- `actual_amount`: The amount of data actually read from the transport layer.
- `optional`: An array containing optional parameters.
- `opt_count`: The length of the array containing optional parameters.

- **Optional parameters:** This function uses optional parameters for some return values.

- `NEAT_TAG_STREAM_ID` (integer): The ID of the stream will be written to this parameter.

- **Return values:** Returns `NEAT_OK` if data was successfully read from the transport layer. Returns `NEAT_ERROR_WOULD_BLOCK` if this call would block. Returns `NEAT_ERROR_MESSAGE_TOO_BIG` if the buffer is not sufficiently large. This is only returned for protocols that are message based, such as UDP, UDP-Lite and SCTP. Returns `NEAT_ERROR_IO` if the connection is reset.

- **Remarks:** This function should only be called from within the `on_readable` callback specified with `neat_set_operations`, as this is the only way to guarantee that the call will not block. NEAT does not permit a blocking read operation. The `actual_amount` value is set to 0 when the remote side has closed the connection.

- **Examples:** None.

- **See also:** `neat_write` and **Optional arguments**

### B.5.9 neat\_write

- **Summary:** Write data to a neat flow. Should only be called from within the `on_writable` callback specified with `neat_set_operations`.

- **Syntax:**

```
1 neat_error_code neat_write(struct neat_ctx *ctx,  
2                          struct neat_flow *flow,  
3                          const unsigned char *buffer,  
4                          uint32_t amount,  
5                          struct neat_tlv optional[],  
6                          unsigned int opt_count);
```

- **Parameters:**

- `ctx`: Pointer to a NEAT context.
- `flow`: Pointer to a NEAT flow.
- `buffer`: Pointer to a buffer containing the data to be written.
- `amount`: The size of the buffer pointed to by `buffer`.
- `optional`: An array containing optional parameters.
- `opt_count`: The length of the array containing optional parameters.

- **Optional parameters:**

- `NEAT_TAG_STREAM_ID` (integer): The ID of the stream the data will be written to.

- **Return values:** Returns `NEAT_OK` if data was successfully written to the transport layer. Returns `NEAT_ERROR_BAD_ARGUMENT` if the specified stream ID is negative. Returns `NEAT_ERROR_OUT_OF_MEMORY` if NEAT is unable to allocate memory. Returns `NEAT_ERROR_WOULD_BLOCK` if this call would block. Returns `NEAT_ERROR_IO` if an I/O operation failed.

- **Remarks:** This function should only be called from within the `on_writable` callback specified with `neat_set_operations`, as this is the only way to guarantee that the call will not block. NEAT does not permit a blocking write operation. Invalid stream IDs are silently ignored.

- **Examples:** None.

- **See also:** `neat_read` and Optional arguments

### B.5.10 neat\_shutdown

- **Summary:** Initiate a graceful shutdown of this flow.
  - The receive buffer can still be read and `on_readable` gets fired like in normal operation.
  - Receiving new data from the peer may fail.
  - All data in the send buffer will be transmitted
  - `neat_write` will fail and `on_writable` will not be called.

- If the peer also has closed the connection, the `on_close` callback gets fired.

- Syntax:

```
1 neat_error_code neat_shutdown(struct neat_ctx *ctx,  
2                             struct neat_flow *flow);
```

- Parameters:

- `ctx`: Pointer to a NEAT context.
- `flow`: Pointer to the NEAT flow to be shut down.

- Return values: Returns `NEAT_OK` if the flow was shut down successfully. Returns `NEAT_ERROR_IO` if NEAT was unable to shut the flow down successfully.

- Remarks: None.

- Examples: None.

- See also: `neat_close`

#### B.5.11 `neat_close`

- Summary: Close this flow and free all associated data. If the peer still has data to send, it cannot be received any more after this call. Data buffered by the NEAT layer which has not given to the network layer yet will be discarded.

- Syntax:

```
1 neat_error_code neat_close(struct neat_ctx *ctx,  
2                             struct neat_flow *flow);
```

- Parameters:

- `ctx`: Pointer to a NEAT context.
- `flow`: Pointer to the NEAT flow to be closed.

- Return values: Returns `NEAT_OK`.

- Remarks: None.

- Examples: None.

- See also: `neat_shutdown`

#### B.5.12 `neat_abort`

- Summary: Abort this flow and free all associated data.

- Syntax:

```
1 neat_error_code neat_abort(struct neat_ctx *ctx,  
2                             struct neat_flow *flow);
```

- Parameters:
  - `ctx`: Pointer to a NEAT context.
  - `flow`: Pointer to the NEAT flow to be aborted.
- Return values: Returns `NEAT_OK`.
- Remarks: Calls `neat_close` internally.
- Examples: None.
- See also: `neat_shutdown` and `neat_close`

### B.5.13 `neat_set_operations`

- Syntax:

```
1 neat_error_code neat_set_operations(  
2     struct neat_ctx *ctx,  
3     struct neat_flow *flow,  
4     struct neat_flow_operations *ops);
```

- Parameters:
  - `ctx`: Pointer to a NEAT context.
  - `flow`: Pointer to a NEAT flow.
  - `ops`: Pointer to a struct that defines the operations/callbacks for this flow.
- Return values: Returns `NEAT_OK`.
- Remarks: struct `neat_flow_operations` is defined as follows:

```
1 struct neat_flow_operations  
2 {  
3     void *userData;  
4  
5     neat_error_code status;  
6     int stream_id;  
7     neat_flow_operations_fx on_connected;  
8     neat_flow_operations_fx on_error;  
9     neat_flow_operations_fx on_readable;  
10    neat_flow_operations_fx on_writable;  
11    neat_flow_operations_fx on_all_written;  
12    neat_flow_operations_fx on_network_status_changed;  
13    neat_flow_operations_fx on_aborted;  
14    neat_flow_operations_fx on_timeout;  
15    neat_flow_operations_fx on_close;  
16    neat_cb_send_failure_t on_send_failure;  
17    neat_cb_flow_slowdown_t on_slowdown;  
18    neat_cb_flow_rate_hint_t on_rate_hint;  
19  
20    struct neat_ctx *ctx;
```

```
21     struct neat_flow *flow;
22 };
```

The information in the `ops` struct will be copied by NEAT.

- **Examples:**

```
1 struct neat_flow_operations ops;
2 ops.on_readable = on_readable;
3 ops.on_writable = on_writable;
4 neat_set_operations(ctx, flow, ops);
```

#### B.5.14 `neat_change_timeout`

- **Summary:** Change the timeout of the flow. Data that is sent may remain un-acked for up to a given number of seconds before the connection is terminated and a timeout is reported to the application.

- **Syntax:**

```
1 neat_error_code
2 neat_change_timeout(struct neat_ctx *ctx, struct neat_flow *flow,
3                    unsigned int seconds);
```

- **Parameters:**

- `ctx`: Pointer to a NEAT context.
- `flow`: Pointer to a NEAT flow.
- `seconds`: The number of seconds after which un-acked data will cause a timeout to be reported.

- **Return values:** Returns `NEAT_OK` if the timeout was successfully changed. Returns `NEAT_ERROR_UNABLE` if attempting to use this function on a system other than Linux, or on flow that is not using TCP. Returns `NEAT_ERROR_BAD_ARGUMENT` if the timeout value is too large or if the specified flow is not opened. Returns `NEAT_ERROR_IO` if NEAT was unable to set the timeout.

- **Remarks:** Only available on Linux for flows using TCP.

- **Examples:** None.

#### B.5.15 `neat_set_primary_dest`

- **Summary:** For multihomed flows, set the primary destination address.

- **Syntax:**

```
1 neat_error_code neat_set_primary_dest(struct neat_ctx *ctx,
2                                     struct neat_flow *flow,
3                                     const char* address);
```

- **Parameters:**



- `ctx`: Pointer to a NEAT context.
- `flow`: Pointer to a NEAT flow.
- `address`: The remote address to use as the primary destination address.
- **Return values:** Returns `NEAT_OK` if the primary destination address was set successfully. Returns `NEAT_ERROR_UNABLE` if the flow is not using SCTP as the transport protocol. Returns `NEAT_ERROR_BAD_ARGUMENT` if the provided address is not a literal IP address.
- **Remarks:** Currently only available for SCTP.
- **Examples:** None.

#### B.5.16 `neat_secure_identity`

- **Summary:** Specify a certificate and key to use for secure connections.
- **Syntax:**

```
1 neat_error_code neat_secure_identity(  
2     struct neat_ctx *ctx,  
3     struct neat_flow *flow,  
4     const char *filename);
```

- **Parameters:**
  - `ctx`: Pointer to a NEAT context.
  - `flow`: Pointer to a NEAT flow.
  - `filename`: Path to the PEM file containing the certificate and key.
- **Return values:** Returns `NEAT_OK`.
- **Remarks:** None.
- **Examples:** None.

#### B.5.17 `neat_set_checksum_coverage`

- **Summary:** Set the checksum coverage for messages sent or received on this flow.
- **Syntax:**

```
1 neat_error_code neat_set_checksum_coverage(  
2     struct neat_ctx *ctx,  
3     struct neat_flow *flow,  
4     unsigned int send_coverage,  
5     unsigned int receive_coverage);
```

- **Parameters:**
  - `ctx`: Pointer to a NEAT context.
  - `flow`: Pointer to a NEAT flow.

- `send_coverage`: UDP-Lite: The number of bytes covered by the checksum when sending messages. UDP: Ignored.
- `receive_coverage`: UDP-Lite: The lowest number of bytes that must be covered by the checksum on a received message. UDP: See below.
- **Return values:** Returns `NEAT_OK` if the checksum coverage was set successfully. Returns `NEAT_ERROR_UNABLE` if the checksum coverage cannot be set, either because the value is invalid, or because the protocol does not support it.
- **Remarks:** Only available for flows using UDP or UDP-Lite. Checksum verification may be enabled/disabled on the receive side for flows using UDP. Specifying a non-zero value for `receive_coverage` will enable it; specifying 0 will disable it.
- **Examples:** None.

#### B.5.18 `neat_set_qos`

- **Summary:** Set the Quality-of-Service class for this flow.
- **Syntax:**

```
1 neat_error_code neat_set_qos(struct neat_ctx *ctx,
2                             struct neat_flow *flow,
3                             uint8_t qos);
```
- **Parameters:**
  - `ctx`: Pointer to a NEAT context.
  - `flow`: Pointer to a NEAT flow.
  - `qos`: The QoS class to use for this flow.
- **Return values:** Returns `NEAT_OK` if the QoS class was set successfully. Returns `NEAT_ERROR_UNABLE` if NEAT was not able to set the requested QoS class.
- **Remarks:** None.
- **Examples:** None.

#### B.5.19 `neat_set_ecn`

- **Summary:** Set the Explicit Congestion Notification value for this flow.
- **Syntax:**

```
1 neat_error_code neat_set_ecn(struct neat_ctx *ctx,
2                              struct neat_flow *flow,
3                              uint8_t ecn);
```
- **Parameters:**
  - `ctx`: Pointer to a NEAT context.
  - `flow`: Pointer to a NEAT flow.

- `ecn`: The ECN value to use for this flow.
- **Return values:** Returns `NEAT_OK` if the QoS class was set successfully. Returns `NEAT_ERROR_UNABLE` if NEAT was not able to set the requested ECN value.
- **Remarks:** None.
- **Examples:** None.

#### B.5.20 `neat_start_event_loop`

- **Summary:** Starts the internal event loop within NEAT.
- **Syntax:**

```
1 neat_error_code neat_start_event_loop(struct neat_ctx *ctx, neat_run_mode run_mode);
```
- **Parameters:**
  - `ctx`: Pointer to a NEAT context.
  - `run_mode`: The mode of which the event loop in NEAT should execute. May be one of either `NEAT_RUN_DEFAULT`, `NEAT_RUN_ONCE`, or `NEAT_RUN_NOWAIT`.
- **Return values:** Returns `NEAT_OK` if the NEAT executed with no error. Returns an error value if the internal event loop in NEAT was stopped due to an error.
- **Remarks:** This function does not return when executed with `NEAT_RUN_DEFAULT`. When executed with `NEAT_RUN_ONCE`, NEAT will poll for I/O, and then block *unless* there are pending callbacks within NEAT that are ready to be processed. These callbacks may be internal. When executed with `NEAT_RUN_NOWAIT`, NEAT will poll for I/O and execute any pending callbacks. If there are no pending callbacks, it returns after polling.
- **Examples:** None.
- **See also:** `neat_stop_event_loop` and `neat_get_backend_fd`

#### B.5.21 `neat_stop_event_loop`

- **Summary:** Stops the internal NEAT event loop.
- **Syntax:**

```
1 int neat_stop_event_loop(struct neat_ctx *ctx);
```
- **Parameters:**
  - `ctx`: Pointer to a NEAT context.
- **Return values:** None.
- **Remarks:** Once called, no further events will be processed and no callbacks will be called until `neat_start_event_loop` is called again.
- **Examples:** None.
- **See also:** `neat_start_event_loop`

### B.5.22 `neat_get_backend_fd`

- Syntax:

```
1 int neat_get_backend_fd(struct neat_ctx *ctx);
```

- Parameters:

- `ctx`: Pointer to a NEAT context.

- Return values: Returns the file descriptor of the event loop used internally by NEAT. May be polled to check for any new events.
- Remarks: Note that embedding this event loop inside another event loop may not be supported on all systems.
- Examples: None.
- See also: `neat_start_event_loop`

### B.5.23 `neat_get_backend_timeout`

- Summary: Return the timeout that should be used when polling the backend file descriptor.

- Syntax:

```
1 int neat_get_backend_timeout(struct neat_ctx *ctx);
```

- Parameters:

- `ctx`: Pointer to a NEAT context.

- Return values: Returns the number of milliseconds on which a poll operation may at most be blocked on the backend file descriptor from libuv before the NEAT event loop should be executed again to take care of timer events within NEAT.
- Remarks: The `client_http_run_once` example demonstrates the use of this function.
- Examples: None.
- See also: `neat_get_backend_fd`

### B.5.24 `neat_get_event_loop`

- Summary: Return the internal NEAT event loop pointer.

- Syntax:

```
1 uv_loop_t neat_get_event_loop(struct neat_ctx *ctx);
```

- Parameters:

- `ctx`: Pointer to a NEAT context.

- Return values: Returns the the event loop used internally by NEAT.
- Examples: None.
- See also: `neat_start_event_loop`

### B.5.25 neat\_get\_stats

- Summary: Return statistics from a NEAT context.

- Syntax:

```
1 neat_error_code neat_get_stats(  
2     struct neat_ctx *ctx,  
3     char **json_stats);
```

- Parameters:

- ctx: Pointer to a NEAT context.
- json\_stats: Pointer to an address where address of the statistics may be written.

- Return values: Returns NEAT\_OK.

- Remarks: The statistics is output in JSON format. The caller is responsible for freeing the buffer containing the statistics.

- Examples: None.

### B.5.26 neat\_getlpasswd

- Summary: Obtains the local or peer addresses of a flow.

- Syntax:

```
1 int neat_getlpasswd(struct neat_ctx* ctx,  
2                   struct neat_flow* flow,  
3                   struct sockaddr** addrs,  
4                   const int local)
```

- Parameters:

- ctx: Pointer to a NEAT context.
- flow: Pointer to a NEAT flow.
- addrs: Pointer to variable for storing pointer to addresses to.
- local: Set to non-zero value for obtaining local addresses, set to 0 to obtain peer addresses.

- Return values: On success, neat\_getlpasswd() returns the number of addresses (local or remote). In case of having obtained at least one address, a pointer to a newly allocated memory area with the addresses will be stored into addrs. This memory area needs to be freed after usage.

- Examples:

```
1 struct struct sockaddr* addrs;  
2 int n = neat_getlpasswd(ctx, flow, &addrs, 1);  
3 if(n > 0) {  
4     struct sockaddr* a = addrs;  
5     for(int i = 0; i < n; i++) {
```

```
6     switch(a->sa_family) {
7         case AF_INET:
8             printf("Address %d/%d: IPv4\n", i, n);
9             a = (struct sockaddr*)((long)a + (long)sizeof(sockaddr_in));
10            break;
11           case AF_INET6:
12               printf("Address %d/%d: IPv6\n", i, n);
13               a = (struct sockaddr*)((long)a + (long)sizeof(sockaddr_in6));
14           default:
15               assert(false);
16               break;
17       }
18   }
19   free(addr);
20 }
```

### B.5.27 neat\_log\_level

- Summary: Set the log-level of the NEAT library.
- Syntax:

```
1 void neat_log_level(struct neat_ctx *ctx,
2                     uint8_t level)
```

- Parameters:

- ctx: Pointer to a NEAT context.
- level: Log level of the log entry

- \* NEAT\_LOG\_OFF
- \* NEAT\_LOG\_ERROR
- \* NEAT\_LOG\_WARNING
- \* NEAT\_LOG\_DEBUG

- Return values: None.
- Examples:

```
1 neat_log_level(ctx, NEAT_LOG_ERROR);
```

- See also: neat\_log\_file

### B.5.28 neat\_log\_file

- Summary: Sets the name of the log file.
- Syntax:

```
1 uint8_t neat_log_file(struct neat_ctx *ctx,
2                       const char* file_name)
```

- **Parameters:**
  - `ctx`: Pointer to a NEAT context.
  - `file_name`: Name of the NEAT logfile. If set to `NULL`, NEAT writes the log output to `stderr`.
- **Return values:** `RETVAL_SUCCESS`: **success**; `RETVAL_FAILURE`: **failure**
- **Examples:**

```
1 neat_log_file(ctx, "disaster.log");
```
- **See also:** `neat_log_level`

## C Main changes with respect to Deliverable D2.2

The following changes log summarises the *main* differences between D2.3 and D2.2 (which this document updates). The list is not exhaustive, and several smaller changes in e.g. the description of NEAT Components are omitted for brevity.

- Section 1.5 was added, providing an overview of connection teardown.
- Section 2.1 was updated with minor modifications to reflect the latest changes in the code and online tutorial.
- Section 2.3: a presentation of Python bindings for the NEAT User API (developed after D2.2 was finished) was added.
- Section 3.2.2 was expanded with new text and figures on the RawRTC library and its integration in NEAT.
- Section 3.2.3 was renamed as “Local Flow Priority” and now includes the case of SCTP with transparent flow mapping.
- Section 3.2.4: for clarity, Table 1 was added showing supported transport protocols and NEAT security in the Core Transport System prototype.
- Section 3.4: the notion of NEAT PropertyArrays was added. Also, new text, figures and tables were added for more clarity, e.g.:
  - Table 3 was added showing REST API endpoints for external access to the Policy Manager.
  - Table 4 was added showing Policy object keys.
  - Figure 12 was added to illustrate the Policy Manager’s inputs and outputs.
- Appendix B was updated to reflect the latest changes in the NEAT API reference, in accordance with the online API reference.



## **D Paper: *NEAT: A Platform- and Protocol-Independent Internet Transport API***

The following research paper [\[29\]](#) has been produced by project participants.

**Preprint (accepted version)**

Published in IEEE Communications Magazine  
<http://www.comsoc.org/commag/>

**Please cite as:**

N. Khademi, D. Ros, M. Welzl, Z. Bozakov, A. Brunstrom, G. Fairhurst, K.-J. Grinnemo, D. Hayes, P. Hurtig, T. Jones, S. Mangiante, M. Tüxen, and F. Weinrank. "NEAT: A Platform- and Protocol-Independent Internet Transport API". *IEEE Communications Magazine*, vol. 55, no. 6, pp. 46–54, 2017.

Digital Object Identifier: <https://doi.org/10.1109/MCOM.2017.1601052>

© 2017 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

# NEAT: A Platform- and Protocol-Independent Internet Transport API

Naeem Khademi, David Ros, Michael Welzl, Zdravko Bozakov, Anna Brunstrom, Gorry Fairhurst, Karl-Johan Grinnemo, David Hayes, Per Hurtig, Tom Jones, Simone Mangiante, Michael Tüxen, and Felix Weinrank

**Abstract**—The sockets Applications Programming Interface (API) has become the standard way that applications access the transport services offered by the Internet Protocol stack. This paper presents NEAT, a user-space library that can provide an alternate transport API. NEAT allows applications to request the service they need using a new design that is agnostic to the specific choice of transport protocol underneath. This not only allows applications to take advantage of common protocol machinery, but also eases introduction of new network mechanisms and transport protocols. The paper describes the components of the NEAT library and illustrates the important benefits that can be gained from this new approach. NEAT is a software platform for developing advanced network applications that was designed in accordance with the standardization efforts on Transport Services (TAPS) in the Internet Engineering Task Force (IETF), but its features exceed the envisioned functionality of a TAPS system.

## I. INTRODUCTION

For more than three decades, the Internet’s transport layer has essentially supported just two protocols and the original design of the sockets API offered only two Transport Services to applications. One service provided stream-oriented in-order reliable delivery, manifested in TCP, and the other a message-based unordered unreliable delivery, manifested in UDP.

Today, more than three decades later, these are the only two transport protocols commonly offered by operating systems to applications. UDP-based applications are used for a wide variety of datagram services from service discovery to interactive multimedia, while TCP became the dominant protocol for Internet services from web browsing to file sharing and video content delivery. While their success has often been attributed to the robustness of these protocols, during the last decades new service requirements have emerged that are beyond what TCP can deliver or UDP can offer—examples include: an interactive multimedia application may prefer to prioritize low latency over strictly reliable delivery of data, but could use partially-reliable delivery to improve quality while ensuring timeliness, or an application may be designed to take

advantage of multihoming when this is available. UDP has also emerged as a substrate upon which user-space transport protocols are being developed—many customized for specific applications (e.g., the QUIC protocol), where much effort can be expended re-implementing common transport functions.

A handful of protocols have been proposed to provide Transport Services beyond those of TCP and UDP; most notably, SCTP, DCCP and UDP-Lite. However none of these have seen widespread use or universal deployment. The reason behind this is often attributed to *ossification* of the Internet’s transport layer, where further evolution has become close to impossible. This has two major aspects:

- **Inflexibility of the current socket API:** Application programmers need to specify *transport protocol-specific* configurations to request a desired service. This binding to protocols inevitably requires programmers to recode their applications to take advantage of any new transport protocol. It also introduces complexity when there is a need to customize for different network scenarios, and choose appropriate transport protocol-specific parameters.
- **Deployment vicious circle:** New protocols and mechanisms cannot be expected to work in unmodified networks. Some equipment may need to be reconfigured, updated or replaced to deploy a new protocol. Developers seeking to use new protocols simply find they cannot be relied upon to work across the Internet. Because the current socket API requires application developers to specifically choose a certain protocol, they therefore tend to avoid using a protocol other than TCP or UDP, knowing that any others are likely to be unsuccessful for many network paths. This chicken-and-egg situation has made it hard for unused transport protocols to become deployed in the Internet—even if they would provide a better service to some applications.

In this paper, we introduce the *NEAT Library*. This is a software library built above the socket API to provide networking applications with a new API offering platform- and protocol-independent access to Transport Services. NEAT is, to the best of our knowledge, the first prototype implementation of IETF standardization efforts on Transport Services (TAPS), which we will discuss in Section V. NEAT and its related standardization efforts in TAPS can re-enable the evolution of the Internet’s transport layer because they break the deployment vicious circle; NEAT’s flexible, customizable API makes it easy to define and use novel services on top of the socket API, seamlessly leveraging new transport protocols

Naeem Khademi and Michael Welzl are with the Department of Informatics, University of Oslo, Norway. E-mail: {naeemk, michawe}@ifi.uio.no.

David Ros and David Hayes are with Simula Research Laboratory, Norway. E-mail: {dros, davidh}@simula.no.

Zdravko Bozakov and Simone Mangiante are with Dell EMC, Ireland. E-mail: {Zdravko.Bozakov, Simone.Mangiante}@dell.com.

Anna Brunstrom, Karl-Johan Grinnemo and Per Hurtig are with Karlstad University, Sweden. E-mail: {anna.brunstrom, karl-johan.grinnemo, per.hurtig}@kau.se.

Gorry Fairhurst and Tom Jones are with the University of Aberdeen, Aberdeen, United Kingdom. E-mail: {tom, gorry}@erg.abdn.ac.uk.

Michael Tüxen and Felix Weinrank are with Münster University of Applied Sciences, Germany. E-mail: {tuexen, weinrank}@fh-muenster.de.

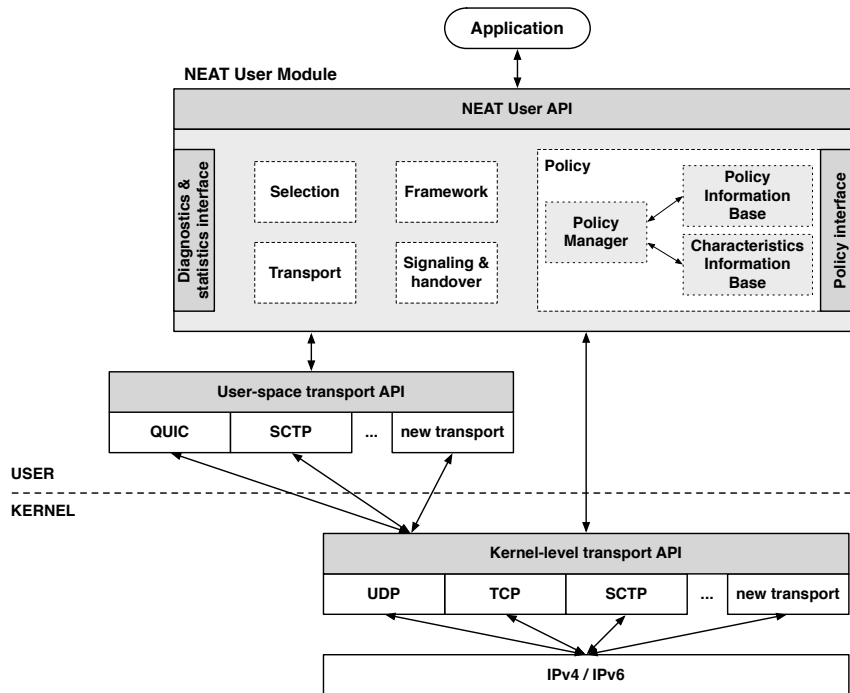


Fig. 1. The architecture of the NEAT System.

as they become available. This, in turn, may create a shift in the traffic pattern seen by vendors and administrators of middleboxes that could at some point lead them to support such traffic. Section IV presents several examples of benefits that NEAT offers to applications.

## II. BACKGROUND

TCP and UDP are a part of the kernel of almost all operating systems and are also supported by nearly all middleboxes. During its lifetime, TCP has been substantially improved. However, the evolution of TCP has had to deal with constraints. Changes to packet format have to consider that middleboxes might block or limit communication. Therefore, a fallback mechanism to the old packet format has become a part of such protocol extensions. New protocol mechanisms (e.g., congestion control or loss recovery mechanisms) mostly focus on single-sided changes to allow faster deployment—but the speed of deployment is still limited by the software development cycle of operating systems.

In addition, having a feature available on an operating system does not imply that it is made available to an application running with user privileges; new features are often disabled by default and turning them on requires special privileges since it has host-wide consequences.

Because UDP provides only minimal services (port numbers and a checksum), it is possible to use it as substrate to implement transport protocols on top of it to introduce features; this approach has become increasingly common. This leads to every UDP-based application to some extent needing to implement the same core set of functions [1]. However, it

also leads to per-application protocol stacks, where transport protocols cannot easily be moved between applications (and making this possible is often not in the interest of the application developer). Developing an efficient transport protocol is a difficult task which requires a number of features to be re-implemented again and again. UDP-based transport protocols have also done nothing to fix the general architectural problem: the socket API's protocol binding remains, typically with a choice between only TCP and UDP.

Specific applications can require services not provided by TCP. One example is the transport of signaling messages in telephony signaling networks. This is used to transfer mostly small messages and requires a high level of fault tolerance. When a protocol stack for this application was developed, a new transport protocol, SCTP [2], was created to fulfill these specific requirements. It was possible to deploy SCTP in these networks because there were no middleboxes, and kernel implementations for the operating systems used in telephony signaling networks were developed.

Currently, the IETF and the World Wide Web Consortium (W3C) are developing WebRTC, a technology for real-time multimedia communication directly between web browsers. Non-media communication using SCTP is also supported; to facilitate deployment across arbitrary Internet paths, SCTP runs over UDP. Google has developed QUIC, a UDP-based transport protocol with features including fast connection setup, cross-layer optimized security, and a modern congestion control and loss recovery mechanism. If QUIC fails to traverse a middlebox, the web browser can fall back to using TCP.

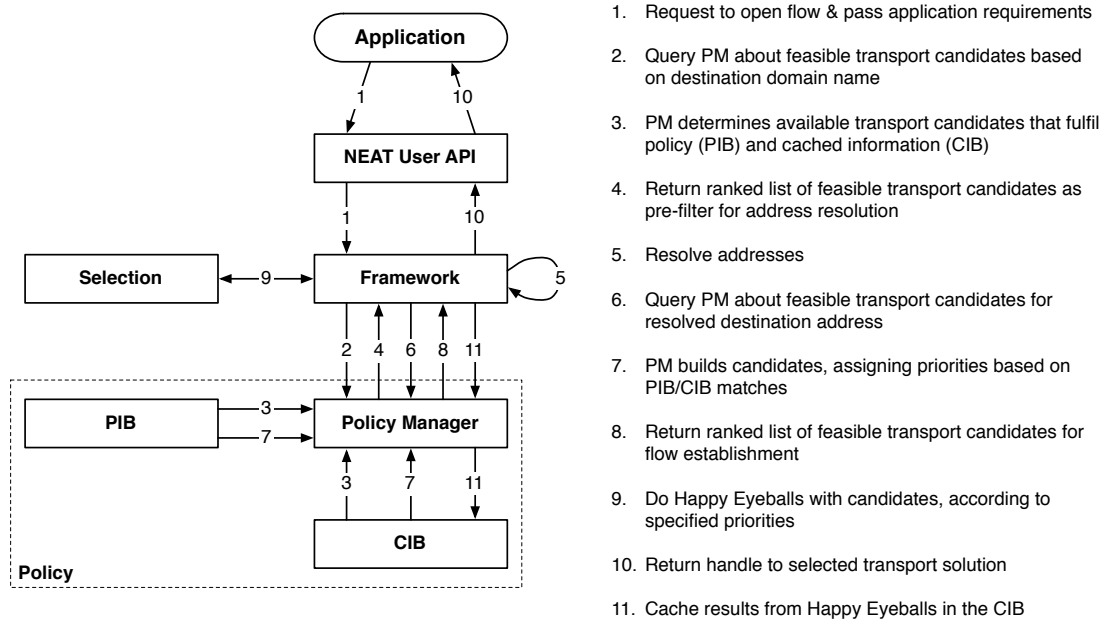


Fig. 2. Simplified workflow showing how NEAT components interact when opening a flow.

By moving a transport protocol from the operating system to the application (e.g., WebRTC and QUIC integrated in web browsers), the release cycle can be substantially shortened, the implementation becomes independent from the operating system, and the protocol can be tailored to the specific application. Using UDP encapsulation is the only option to not using TCP and be able to traverse middleboxes.

This enables a larger variety of transport protocols to co-exist and change over time, but does not help with the issue of per-application protocol stacks mentioned before. An application programmer has to add complexity to benefit from advanced features in their application; this requires utilizing different APIs, figuring out which protocols are supported by the remote end-points, selecting protocol mechanisms, and providing fallback mechanisms when these happen to not work across the current network path. None of these are specific to a particular transport protocol, but are related to the need for the programmer to work with a variety of transport protocols. This general problem could be addressed by defining a new transport system, as outlined next.

### III. RE-ENABLING EVOLUTION: INTRODUCING NEAT

As discussed above, using transport services beyond TCP and UDP today puts a high burden on the application developer. The NEAT Library addresses this problem by providing application developers with one enhanced API that is transport-protocol independent, with the library providing support for selecting the best available transport option *at run-time* and handling fallback between transport protocols as needed. Running as a user-space library, NEAT can make use of transports running both in user space and in the kernel, all transparent to the applications. Protocols like SCTP are

already supported over many paths, but they cannot be easily used by application programmers unless they are supported over *all* paths. NEAT changes this by placing functions such as selecting a transport and handling fallback below the API. NEAT allows such functions to evolve with the network, rather than be bound to specific applications.

Figure 1 provides a schematic view of the NEAT architecture. Applications employ the NEAT User API to access transport services. This API is located in the NEAT User Module, which is the core of NEAT and comprises components that together deliver services tailored to application requirements at run-time. The components in the module are grouped in five categories: Framework, Policy, Selection, Transport, and Signaling & Handover.

Framework components provide basic functionality required to use NEAT. They define the structure of the NEAT User API and implement core library mechanisms. Applications provide information about their requirements for a desired transport service via this API.

Policy components comprise the Policy Information Base (PIB), the Characteristics Information Base (CIB), and the Policy Manager (PM). The function of the PM is to generate a ranked list of connection candidates that fulfill the application requirements while taking system and network constraints into account and adhering to configured policies. All policy components operate on so-called NEAT Properties, which express requirements and characteristics throughout the NEAT System. Each property is a key-value tuple with additional metadata indicating the priority (mandatory or optional) and weight of the associated attribute.

Policies and profiles—stored in the PIB—extend and modify the property set associated with each connection candidate.

In addition, the CIB repository maintains information about available interfaces, supported protocols towards previously-accessed destination endpoints, network properties and current/previous connections between endpoints. The content of the CIB is continuously updated by local and external CIB sources.

Selection components choose an appropriate transport solution. The additional information provided by the NEAT User API enables the NEAT Library to move beyond the constraints of the traditional socket API, making the stack aware of what is actually *desired* or *required* by the application. On the basis of both the information provided by the NEAT User API and the PM, candidate transport solutions are identified. The candidate solutions are then tested by the Selection components, and the one deemed most appropriate is then used.

Transport components are responsible for providing functions to instantiate a transport service for a particular traffic flow. They provide a set of transport protocols and other necessary components to realize a transport service. While the choice of transport protocols is handled by the Selection components, the Transport components are responsible for configuring and managing the selected transport protocols.

Signaling & Handover components can provide advisory signaling to complement the functions of the Transport components. This could include communication with middleboxes, support for handover, failover and other mechanisms.

Figure 2 illustrates a simplified workflow, showing how the NEAT components interact when an application initiates a new flow. As follows from the above description, NEAT has an evolvable architecture that opens up for the introduction of new transport services and can enable interaction with network devices to improve such services. NEAT also enables the incremental introduction of new transport protocols, both in the kernel and in user space, as the API is independent from the underlying transport protocol.

#### IV. BENEFITS OF NEAT

Next, we present four examples of key benefits of using the NEAT Library. First, NEAT provides an API that is simple to use. This allows existing applications to be easily ported to the NEAT Library, simplifying network communication and reducing code complexity.

NEAT also provides automatic fallback using a *Happy Eyeballs* (HE) mechanism. HE is a generic term for algorithms that test for end-to-end support of a protocol X simply by trying to *use* X, then falling back to a default choice Y known to work if X is found to not work (e.g., after a suitable timeout). This added functionality is lightweight and has negligible cost compared to other communication tasks. It allows applications to take advantage of the best available transport solution and in turn enables transport innovation (e.g., applications do not need to be recoded to use a new transport feature or protocol that becomes available).

NEAT not only facilitates evolution of the transport protocols and introduction of new transport mechanisms, it can also help enable innovation at the network layer. The higher-level of abstraction offered by the NEAT User API eases

the path to utilizing Quality of Service (QoS) support for UDP-based applications, and could be used to access other network services should they become available (e.g., selection of the most cost-effective or secure path utilizing IPv6 provisioning-domain information). Applications and networks can also leverage the flexible control provided by the Policy components, for example to provide a generic interface for exchanging information between external SDN controllers and NEAT-enabled applications.

##### A. Porting applications to NEAT

The NEAT User API offers a uniform way to access networking functionality, independent from the underlying network protocol or operating system. Many common network programming tasks like address resolution, buffer management, encryption, connection establishment and handling are built into the NEAT Library and can be used by any application that uses NEAT.

Developers write applications using the asynchronous and non-blocking NEAT User API, implemented using the libuv [3] library which provides asynchronous I/O across multiple-platforms.

As shown in Listing 1, users can request the services that they expect from the network (e.g. low latency, reliable delivery, a specific TCP congestion control algorithm) by providing an optional set of properties to control the behavior of the library.

```
static neat_error_code
on_connected(struct neat_flow_operations *ops)
{
    // set callbacks to write and read data
    ops->on_writable = on_writable;
    ops->on_all_written = on_all_written;
    ops->on_readable = on_readable;
    neat_set_operations(ops->ctx, ops->flow, ops);
    return NEAT_OK;
}

int
main(int argc, char *argv[])
{
    // initialization of basic NEAT structures
    struct neat_ctx *ctx;
    struct neat_flow *flow;
    struct neat_flow_operations ops;
    ctx = neat_init_ctx();
    flow = neat_new_flow(ctx);
    memset(&ops, 0, sizeof(ops));

    // callback when connection is established
    ops.on_connected = on_connected;
    neat_set_operations(ctx, flow, &ops);

    // optional user requirements in JSON format
    static char *properties = "{\"transport\":{\"SCTP\", \"TCP\"}}";
    neat_set_property(ctx, flow, properties);

    // connect
    if (neat_open(ctx, flow, "127.0.0.1", 5000, NULL, 0) {
        fprintf(stderr, "neat_open failed\n");
        return EXIT_FAILURE;
    }

    // start libuv loop
    neat_start_event_loop(ctx, NEAT_RUN_DEFAULT);

    neat_free_ctx(ctx);

    return EXIT_SUCCESS;
}
```

Listing 1. Code example from a simple client using the NEAT API.

The NEAT Library then uses a set of internal components to establish a connection over the network. To make an

appropriate selection, the Policy Manager maps user properties to policies and computes a set of candidate transports that can satisfy the request. NEAT also can utilize policy information directly set by the user, system administrator or developer.

Connections to a peer endpoint are made by creating a new flow, which is a bidirectional link between two endpoints similar to a socket in the traditional Berkeley Socket API but not strictly tied to an underlying transport protocol.

The NEAT API executes callbacks in the application when an event from the underlying transport happens, creating a more natural and less error-prone way of network programming than with the traditional socket API. The three most important callbacks in the NEAT API are `on_connected`, called once the flow has connected to a remote endpoint; `on_readable` and `on_writable`, called once data may be written to or read from the flow.

Our experience with NEAT shows a reduction of the code size by  $\approx 20\%$  for each application, as the library streamlines a number of connection establishment steps. For example, the single function call `neat_open` requests name resolution and all other functions required before communication can start, hiding complex boilerplate code. Ported applications remain fully interoperable with regular TCP/IP-based implementations, while being able to take advantage of NEAT functionalities. Besides, they can benefit from support for alternative transports, when available, relieving programmers from dealing with fallbacks between protocols. Finally, a traditional socket-based shim layer has been implemented on top of NEAT to allow legacy applications to make use of NEAT functionalities through policies without requiring direct porting to the NEAT API.

### B. Happy Eyeballs: A Lightweight Transport Selection Mechanism

Selection components employ a HE mechanism to enable a source host to determine whether a transport protocol is supported along the current network path. This allows applications to benefit from advances in transports that may be only partially deployed in the Internet. The HE mechanism used by NEAT is similar to that introduced to facilitate IPv6 adoption [4], but works at the transport layer to select one of a set of connection-oriented transport solutions. The Selection components receive a ranked list of potential candidates generated by the PM, where a higher ranking indicates a better match with application and policy requirements. The HE mechanism then concurrently tries each transport solution from the list, delaying initiation of lower-priority transport solutions.

Figure 3 shows the HE mechanism in a scenario where the best transport to the destination is unknown and current policy dictates that the HE process is used to select between TCP and SCTP, but preferring SCTP. The initiation of the TCP connection is delayed for a time interval governed by policy, specifying a difference in priority between candidate protocols. If the SCTP connection does not complete within the time interval, a TCP connection is also started. The first transport to complete a connection is selected and becomes the transport

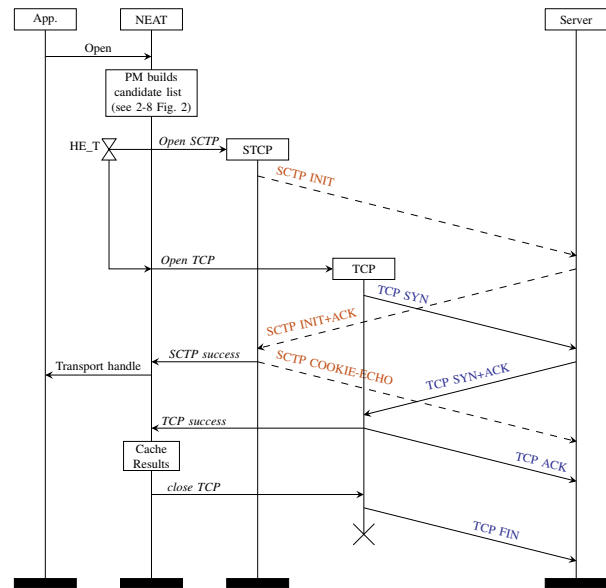


Fig. 3. Message Sequence Chart (MSC) illustrating the NEAT Happy Eyeballs (HE) transport selection process when selecting between TCP and SCTP, SCTP preferred.

of choice. Once connectivity is established, other methods are abandoned, and their connections closed.

To avoid wasting network resources by routinely attempting concurrent connections, HE instructs the Policy components to cache the outcome of each selection result in the CIB for a configurable amount of time. After expiry of the time, the selection is removed from the cache, re-enabling HE.

Consider the scenario in Figure 3. Attempting selection when there is no existing cache entry requires extra resources, potentially resulting in opening connections for each candidate transport protocol. In this example, SCTP completes first and the TCP connection is closed having sent no data. With typical web traffic and worst-case packet sizes, byte overhead is as small as  $\approx 1\%$ . For a cache hit rate of 80%, this reduces further to  $\approx 0.2\%$ . A detailed evaluation of the impact of HE in terms of memory and CPU utilization can be found in [5], where it is shown that CPU costs are relatively small (especially when considering the cost of TLS encryption), and that HE has only a minor impact on memory consumption.

### C. Deployable QoS with NEAT

Network QoS is often used for traffic engineering, but few applications have managed to exploit this technology beyond a controlled network environment. One major obstacle is the lack of a consistent high-level API.

There have been attempts to add methods that directly associate QoS with IP traffic (e.g., [6], [7]), but they have seen little to no adoption. A key challenge is how to express the service requirements, while still enabling policy to influence choice and providing flexibility when the network is unable to directly satisfy the requirements.

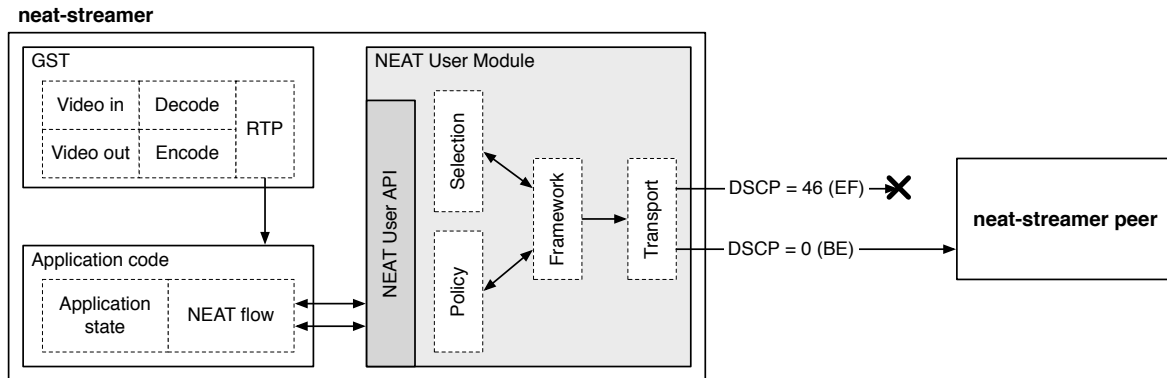


Fig. 4. Example of neat-streamer using QoS fallback with NEAT. The application sets up the media pipeline and uses NEAT to transfer data across the network, according to the requested service. The NEAT Library could try to send UDP datagrams with a DSCP set to High Priority Expedited Forwarding (EF). A timer triggers the NEAT Library to query the application status, which then reveals the application failed to use this DSCP, so NEAT can now try the next DSCP value, Default (BE). When the timer again triggers, the application reports success and this code point continues to be used.

The NEAT API can allow applications to specify QoS requirements. This can, for example, utilize policy information to drive an appropriate Differentiated Service Code Point (DSCP). The finally chosen DSCP can be based on both static policy and dynamic information collected from connections using NEAT.

The NEAT fallback mechanism can be used with any datagram services to enable the NEAT Library to select between a list of candidate datagram transports, network encapsulations and interfaces. This can assist an application to robustly find desirable connection parameters for any path by transparently falling back to alternative services when required (resembling, but different to the NEAT HE function for connection-oriented transports).

Neat-streamer [8] is a demo application that utilizes the NEAT Library for live streaming video over connectionless transports using the GStreamer (GST) media libraries. GST is a pipeline-based media system that supports a wide range of audio and video formats and other functions via a plugin system.

Figure 4 shows the interactions between NEAT and neat-streamer running on a network that drops traffic with certain DSCP values set.

Because neat-streamer uses NEAT, it can indicate the QoS treatment that it requires for each media flow, and the endpoint to which it wishes to stream. NEAT provides the required QoS marking and may determine which transport service to use (e.g., choosing between UDP-Lite, UDP, or use of Traversal Using Relays around NAT, TURN), and whether security functions are required.

NEAT also provides the protocol machinery to update the selected flow parameters should network connectivity problems be reported by the application. A timer triggers a callback function within the application to determine whether the application believes the network is delivering the service it requires (in many cases, only the application is aware of the performance reported by a remote datagram receiver). When an application reports failure it can allow NEAT to use the list of candidates, and potentially other information (e.g., held

within the CIB) to search for alternate parameters.

#### D. SDN Integration

The ability of enabling external sources to query and augment the state of the Policy Manager is a key design choice of the NEAT architecture. As a consequence, NEAT-enabled end-hosts can be seamlessly integrated in centrally controlled environments, such as Software-Defined Networks (SDNs). In such environments, logically centralized controllers aim to maintain a global view of the network and optimize its utilization. To achieve this, controllers ideally require detailed and up-to-date knowledge of available resources, in addition to the requirements and characteristics of deployed applications. Today, controllers rely on time-consuming and error-prone heuristics to infer the association between applications, their requirements, and observed flows.

In this context, the benefit of the NEAT approach is three-fold. Firstly, NEAT applications may inform controllers directly about their particular requirements towards the network. In NEAT, such requirements are defined either explicitly by application developers, or through suitable system policies. This strategy can reduce the need for network controllers to guess how to treat individual flows. Secondly, through the Policy Manager CIB, NEAT enables controllers to supply applications with detailed information about the state of paths available to the host. In the absence of this feedback, metrics such as available bandwidth or latency may need to be inferred individually by each application through measurements. Finally, the controller gains the ability to deploy policies at the host level which influence the transport protocols, interfaces and associated parameters used in NEAT applications.

All mechanisms necessary for exchanging information between the controller and NEAT-enabled applications are implemented in Policy components. Specifically, the Policy Interface is exposed through a REST API, enabling external entities to push information to the PIB and CIB and query their contents. As a result, for each flow request created by a NEAT application, the Policy Manager will utilize the latest policies



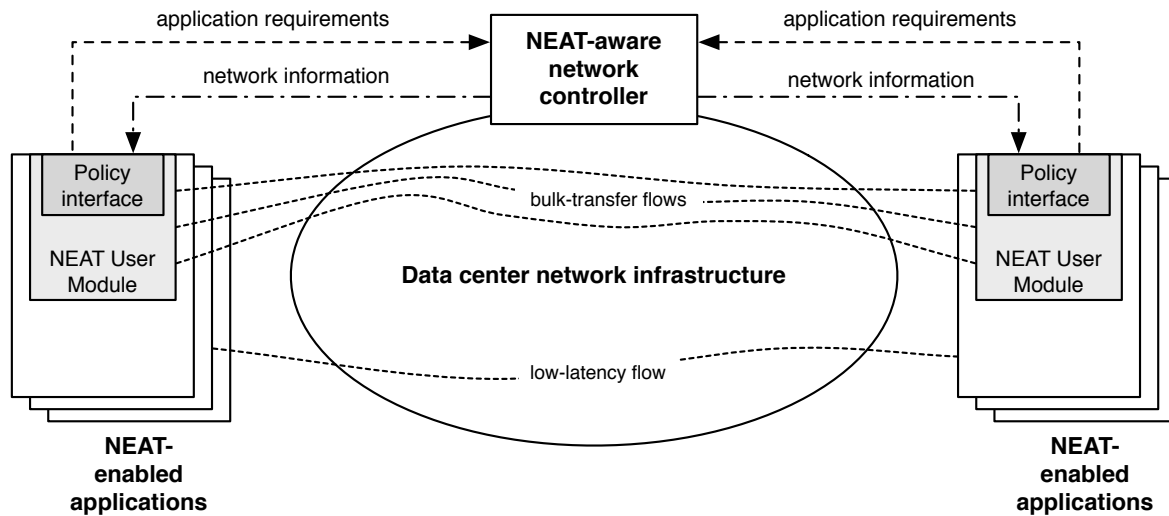


Fig. 5. SDN architecture in which a controller uses the NEAT Library to supply end-hosts with information about the available network resources, and to collect information about the application requirements.

and network attributes supplied by a controller to select the most suitable connection option. Similarly, the API allows the controller to query the CIB to identify the requirements associated with specific application flows or relevant policies configured in the PIB.

To demonstrate the feasibility of the aforementioned controller integration we have implemented a scenario comprised of NEAT-enabled hosts deployed in an OpenFlow SDN network. The aim of the scenario, depicted in Figure 5, is to enable a controller to steer the handling of bulk traffic flows. Each host runs a NEAT-enabled data replication application which provides the estimated flow size as part of the NEAT API call. We used the OpenDaylight framework to implement a controller which monitors the network utilization and calculates a data volume threshold above which flows are considered as bulk flows. We implemented a northbound API and periodically publish a policy to the NEAT end hosts. The policy is triggered when the flow size exceeds the threshold and forces the flows to be tagged with a predefined DSCP marking. As a result, flows affected by the policy are routed through a pre-provisioned network path.

## V. STANDARDIZATION

Recognizing the need for the transport layer (socket) interface to become protocol-independent, the IETF chartered a working group called “Transport Services” (TAPS) in September 2014. A common approach in prior work was to start analysis based on the needs of applications. Instead, TAPS used a methodology that started from a survey of the services offered by available IETF transport protocols [9]. It is currently documenting the primitives and parameters used to access features of a subset of these protocols [10] to form a basis for the design of a protocol-independent API. NEAT developers have been actively contributing to this initiative based on

experience of using the NEAT API, which shares many of the goals behind development of TAPS.

The working group is now shortening the list of transport features. Examples of features include “Specify ECN field” or “Choice between unordered (potentially faster) or ordered delivery of messages”. A recent contribution by NEAT developers [11] recommends against exposing a transport feature in the API when either choosing or configuring it requires knowledge specific to the network path or the operating system, but not the application. A final step will eliminate features specific to a particular protocol that cannot reasonably be implemented using a different protocol—such features contradict the main purpose of TAPS, to be protocol-independent. At the end of this process, this will result in a subset of transport features that end systems supporting TAPS need to provide. NEAT implements all services specified in the current TAPS documents and may therefore be regarded as a prototype implementation of TAPS.

TAPS is also chartered to define experimental support mechanisms, for example to select and engage an appropriate protocol and discover the set of protocols available for a selected service between a given pair of endpoints, to allow the operating system to choose between protocols (e.g., HE and application-level feedback mechanisms). This approach of breaking the binding between applications and transport protocols is an important final step for TAPS.

## VI. CONCLUSION

The service needs of today’s Internet applications range well beyond the basic ones provided by TCP and UDP. Yet, the Internet’s transport layer, as it presents itself to a developer via the socket API, has remained unchanged. This has led to per-application (and per-company) developments in user space, over UDP, such as QUIC for Google Chrome. While these new UDP-based transport protocols have recently

pushed the transport layer into the spotlight, they are also only silo solutions which do nothing to solve the architectural ossification problem: the socket API binds applications to protocols at design time—therefore, transport protocols cannot be replaced without changing applications.

In this paper we presented the NEAT Library, which lets application developers access features of transport protocols in a simple and uniform way. NEAT helps freeing developers from platform or protocol dependencies; they do not have to worry about the specifics of each protocol or operating system; they also do not need to worry about whether a protocol works on a given path. Underneath the NEAT User API, new protocols can seamlessly be inserted, automatically yielding benefits to the application on top. With NEAT's clear layer separation, the Internet's transport layer can finally evolve again.

At the time of writing, prototype code for all component types has been developed for several Unix-like OSs. Besides `neat-streamer`, the NEAT development team has ported example applications to NEAT for early testing, including the `Nghttp2` [12] web server and client, several smaller applications like HTTP/HTTPS clients and performance measurement tools; also, a NEAT-supported Firefox implementation is currently under development by Mozilla. NEAT is an open-source project that welcomes contributions. Source code, documentation and implementation status can be found on GitHub [13].

#### ACKNOWLEDGMENT

The authors would like to thank the anonymous reviewers for their useful remarks.

This work has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No. 644334 (NEAT). The views expressed are solely those of the author(s).

#### REFERENCES

- [1] L. Eggert and G. Fairhurst, "Unicast UDP Usage Guidelines for Application Designers," RFC 5405 (Best Current Practice), Internet Engineering Task Force, Nov. 2008, accessed on February 23, 2017. [Online]. Available: <http://www.ietf.org/rfc/rfc5405.txt>
- [2] R. Stewart, "Stream Control Transmission Protocol," RFC 4960 (Proposed Standard), Internet Engineering Task Force, Sep. 2007, accessed on February 23, 2017. [Online]. Available: <http://www.ietf.org/rfc/rfc4960.txt>
- [3] libuv — Cross-platform Asynchronous I/O. Accessed on February 23, 2017. [Online]. Available: <https://libuv.org/>
- [4] D. Wing and A. Yourtchenko, "Happy Eyeballs: Success with Dual-Stack Hosts," RFC 6555 (Proposed Standard), Internet Engineering Task Force, Apr. 2012, accessed on February 23, 2017. [Online]. Available: <http://www.ietf.org/rfc/rfc6555.txt>
- [5] G. Papastergiou, K.-J. Grinnemo, A. Brunstrom, D. Ros, M. Tüxen, N. Khademi, and P. Hurtig, "On the Cost of Using Happy Eyeballs for Transport Protocol Selection," in *Proceedings of the 2016 Applied Networking Research Workshop (ANRW)*. Berlin: ACM, Jul. 2016, pp. 45–51.
- [6] H. Abbasi, C. Poellabauer, K. Schwan, G. Losik, and R. West, "A Quality-of-service Enhanced Socket API in GNU/Linux," in *4th Real-Time Linux Workshop*, 2002, accessed on February 23, 2017. [Online]. Available: [https://www.osadl.org/fileadmin/events/rtlws-2002/proc/g08\\_abbasi.pdf](https://www.osadl.org/fileadmin/events/rtlws-2002/proc/g08_abbasi.pdf)
- [7] P. Gomes Soares, Y. Yemini, and D. Florissi, "QoSockets: A New Extension to the Sockets API for End-to-end Application QoS Management," *Computer Networks*, vol. 35, no. 1, pp. 57–76, 2001.
- [8] Neat-streamer Video Workload Tool. Accessed on February 23, 2017. [Online]. Available: <https://github.com/uaaerg/neat-streamer>
- [9] G. Fairhurst, B. Trammell, and M. Kühlewind, "Services provided by IETF transport protocols and congestion control mechanisms," Internet Engineering Task Force, Internet-Draft draft-ietf-taps-transport-11, Sep. 2016, work in Progress. Accessed on February 23, 2017. [Online]. Available: <https://tools.ietf.org/html/draft-ietf-taps-transport-11>
- [10] M. Welzl, M. Tüxen, and N. Khademi, "On the Usage of Transport Service Features Provided by IETF Transport Protocols," Internet Engineering Task Force, Internet-Draft draft-ietf-taps-transport-usage-01, Jul. 2016, work in Progress. Accessed on February 23, 2017. [Online]. Available: <https://tools.ietf.org/html/draft-ietf-taps-transport-usage-01>
- [11] S. Gjessing and M. Welzl, "A Minimal Set of Transport Services for TAPS Systems," Internet Engineering Task Force, Internet-Draft draft-gjessing-taps-minset-03, Oct. 2016, work in Progress. Accessed on February 23, 2017. [Online]. Available: <https://tools.ietf.org/html/draft-gjessing-taps-minset-03>
- [12] T. Tsujikawa, `Nghttp2: HTTP/2 C Library`. <https://nghttp2.org/>. Accessed on February 23, 2017.
- [13] NEAT GitHub public repository. Accessed on February 23, 2017. [Online]. Available: <https://github.com/NEAT-project/neat>

## **E Paper: *On the Cost of Using Happy Eyeballs for Transport Protocol Selection***

The following research paper [\[33\]](#) has been produced by project participants.

# On the Cost of Using Happy Eyeballs for Transport Protocol Selection

Giorgos Papastergiou<sup>†</sup>, Karl-Johan Grinnemo<sup>‡</sup>, Anna Brunstrom<sup>‡</sup>, David Ros<sup>†</sup>,  
Michael Tüxen<sup>\*</sup>, Naeem Khademi<sup>\*</sup>, Per Hurtig<sup>‡</sup>

<sup>†</sup>Simula Research Laboratory, <sup>‡</sup>Karlstad University, <sup>\*</sup>Fachhochschule Münster, <sup>\*</sup>University of Oslo

{gpapaste, dros}@simula.no, {karl-johan.grinnemo, anna.brunstrom,  
per.hurtig}@kau.se, tuexen@fh-muenster.de, naeemk@ifi.uio.no

## ABSTRACT

Concerns have been raised in the past several years that introducing new transport protocols on the Internet has become increasingly difficult, not least because there is no agreed-upon way for a source end host to find out if a transport protocol is supported all the way to a destination peer. A solution to a similar problem—finding out support for IPv6—has been proposed and is currently being deployed: the *Happy Eyeballs* (HE) mechanism. HE has also been proposed as an efficient way for an application to select an appropriate transport protocol. Still, there are few, if any, performance evaluations of transport HE. This paper demonstrates that transport HE could indeed be a feasible solution to the transport support problem. The paper evaluates HE between TCP and SCTP using TLS encrypted and unencrypted traffic, and shows that although there is indeed a cost in terms of CPU load to introduce HE, the cost is relatively small, especially in comparison with the cost of using TLS encryption. Moreover, our results suggest that HE has a marginal impact on memory usage. Finally, by introducing caching of previous connection attempts, the additional cost of transport HE could be significantly reduced.

## CCS Concepts

•Networks → Transport protocols; Network performance evaluation;

## Keywords

Transport-protocol selection, Happy Eyeballs, TCP, SCTP, TLS, CPU load, memory usage.

## 1. INTRODUCTION

The deployment of new transport protocols on the Internet is not a trivial task. Several hurdles have to be cleared before a new transport can be used between an arbitrary

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or to publish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

ANRW '16, July 16 2016, Berlin, Germany

© 2016 Copyright held by the owner/author(s). Publication rights licensed to ACM. ISBN 978-1-4503-4443-2/16/07...\$15.00

DOI: <http://dx.doi.org/10.1145/2959424.2959437>

pair of end hosts and see wide adoption. One of the main issues that has to be solved is, how can an end host *know* if a new protocol X is supported along the whole end-to-end path, including the remote host? In the absence of a priori knowledge or explicit signaling, the only way to know whether X works is to *try* it.

Testing a set of candidate protocols can be done serially—e.g., try first with the preferred choice X and, if the attempt fails after a suitable timeout, then fall back to a default alternative Y. Since a connection timeout can introduce a delay of up to tens of seconds, serializing attempts can incur a large latency penalty when the new protocol X is not supported, stalling the application until the subsequent connection trial succeeds.

The *Happy Eyeballs* (HE) mechanism was introduced as a means to facilitate IPv6 adoption [13]. Dual-stack client applications should be encouraged to try setting up connections over IPv6 first, and fall back to using IPv4 if IPv6 connection attempts fail. However, serializing tests for IPv6 and IPv4 connectivity can result in large connection latency. Happy Eyeballs for IPv6 minimizes the cost in delay by *parallelizing* attempts over IPv6 and IPv4.

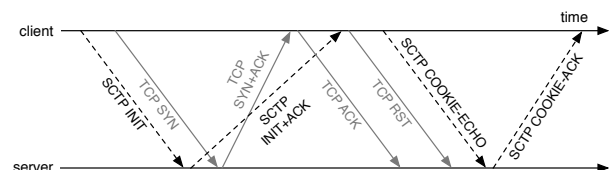


Figure 1: HE for selecting between SCTP and TCP, with SCTP being the preferred choice.

The basic idea behind Happy Eyeballs for IPv6 can be extended to discover support for transport protocols, and in particular to allow an application to use SCTP when it is available end-to-end, else revert to TCP when it is not [11, 12]. Figure 1, adapted from [12], depicts how HE for transport selection may work. An end host simultaneously initiates a TCP connection and an SCTP association; it is assumed that the SYN+ACK arrives before the INIT-ACK. If SCTP (the preferred choice) is supported, the TCP connection is abandoned<sup>1</sup>. Similar to what is done in actual

<sup>1</sup>The figure assumes there is a mechanism to notify the application about the reception of the SCTP INIT-ACK, so the

implementations of HE for IPv6 [10], a small delay may be introduced to give an advantage to SCTP over TCP [14]; else, the host can just pick the protocol for which a response (TCP SYN+ACK, SCTP INIT-ACK) arrives first.

This paper focuses on the performance penalties introduced by HE for transports. In particular, we study the impact that a HE mechanism for selecting between TCP and SCTP may have on CPU load and memory usage at a destination end-host. Our results provide empirical arguments in favor of using such a mechanism for transport selection.

The rest of the paper is organized as follows. Section 2 provides some background and motivation for introducing a transport HE mechanism. As follows from this discussion, two major concerns about transport “happy-eyeballing” are higher CPU load and memory usage. In Section 3 we assess these concerns by experimentally evaluating HE between TCP and SCTP. Finally, Section 4 concludes the paper.

## 2. BACKGROUND

Some thirty years back, there were two Internet transport protocols to choose from: TCP and UDP. Since these two protocols basically represent the opposites in the services they provide—TCP provides a reliable, in-order, byte-stream oriented delivery service and UDP an unreliable, unordered message delivery—the selection between these two protocols from an application viewpoint was mostly straightforward. Furthermore, it could be expected that all hosts supported both TCP and UDP, and there were no middleboxes altering or blocking the traffic before it reached its final destination.

Today, the Internet looks rather different. The number of standard transport protocols and their options (and the different services they may provide) has increased, making the selection of a suitable transport less straightforward. New transports may allow to provide improved services to applications, but middleboxes such as firewalls, NATs and load balancers have become an integral part of the Internet, and there is a great diversity in how they are configured and deployed; it cannot be assumed that any transport or transport option can safely make it from sender to receiver.

As mentioned in Section 1, although HE was primarily introduced as a way to promote the use of IPv6, it has also been proposed as a way for an application to efficiently select transports [11, 12]. Wing and Yourtchenko [11] provide recommendations for HTTP clients on how to seamlessly migrate from TCP to SCTP without any adverse impact on the user experience. Moreover, they propose a way to combine an IPv6/IPv4 HE with a TCP/SCTP HE for a web browser running on a dual-stack machine [12]. Also worth mentioning in this context is the work carried out by the Transport Services (TAPS) working group of the IETF [8]. One of this working group’s planned documents should “[...] explain how to select and engage an appropriate protocol and how to discover which protocols are available for the selected service between a given pair of end points [...]”, something which will likely require HE between transport solutions.

Still, a HE transport-selection mechanism does raise questions about increased CPU usage and memory consumption.

application can then abort the TCP connection by sending a TCP Reset. Without such mechanism, the TCP connection can only be aborted after the full four-way handshake of SCTP is completed.

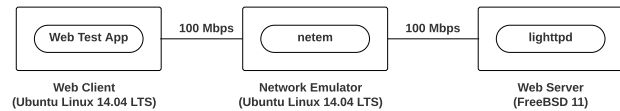


Figure 2: Experiment setup.

When HE is used, a single connection request from the application might result in several concurrent transport connection requests, i.e., not just one connection request at a time as is the case when HE is not used. Hence, the use of HE could result in an increase in both CPU and memory usage. Baker [3] provides recommendations on how to evaluate IPv4/IPv6 HE, however, metrics like CPU load and memory usage are not considered in [3]. Note that neither of these are key metrics for IP HE, whereas they are for transport HE—this is so because transport connection setup means creating state in end points. There have been a few discussions of the performance of HE for IPv6/IPv4 [1, 2, 4, 7] but, to the best of our knowledge, this paper is the first to focus on performance aspects of transport-layer HE, in terms of CPU load and memory usage.

A HE transport-selection mechanism also raises questions about increased use of network resources, a key issue for the scalability of HE. For instance, the aforementioned HE proposal by Wing et al. [12] transmits four packets for every application connection request. Still, as already pointed out by Wing et al. [11, 13, 14], HE network resource usage should be mitigated by the use of caching.

## 3. EVALUATION

This section evaluates HE between TCP and SCTP in terms of CPU load and memory usage. The section begins with a description of the experiment setup and the studied test scenarios. The remainder of the section presents and comments on the results from the execution of these scenarios.

### 3.1 Experiment Setup

In our experiment, we modeled a single wide-area network path to an upstream Web server. The laboratory network used in our experiment is shown in Figure 2. The three machines in the experiment were of type: Dell Optiplex 9020 with 3.60 GHz Intel Core i7-4790 (quad core) processors. The Web Client and Network Emulator machines ran Ubuntu Linux 14.04 LTS with kernel 3.13.0, and the Web Server machine ran FreeBSD 11 (revision r294499). All machines used the default network kernel settings, except those listed in Table 1. These changes assured that the testbed could properly support the connection rates considered in this work, and disabled all SCTP features that were not needed in the experiments. The Network Emulator machine used *netem* to emulate a propagation delay of 20 ms.

The Web Client hosted a custom-designed Web traffic generator in which two modified versions of the *httperf* [6] web traffic generator (one that supports TCP and one that supports SCTP) were combined to implement the studied test scenarios. HTTP/1.0 with the Keep-Alive option enabled was used in both *httperf* programs. The FreeBSD server hosted a *lighttpd* [9] server, modified to listen for both TCP and SCTP HTTP/1.0 unencrypted and TLS-encrypted requests. The *lighttpd* server was also modified

Table 1: Kernel Settings

Web Server	Settings
net.inet.tcp.syncache.hashsize	2048
kern.ipc.somaxconn	4096
net.inet.sctp.pr_enable	0
net.inet.sctp.ecn_enable	0
net.inet.sctp.outgoing_streams	1
net.inet.sctp.incoming_streams	1
net.inet.sctp.asconf_enable	0
net.inet.sctp.auth_enable	0
net.inet.sctp.reconfig_enable	0
net.inet.sctp.nrsack_enable	0
net.inet.sctp.pktdrop_enable	0
Web Client	Settings
net.ipv4.ip_local_port_range	10000 61000
net.ipv4.tcp_tw_recycle	1
Network Emulator	Settings
net.ipv4.ip_forward	1

Table 2: lighttpd Settings

Configuration Parameter	Settings
server.network-backend	writenv
server.event-handler	kqueue
server.max-fds	4096
server.max-connections	2048
server.max-worker	7
ssl.use-sslv2	disable
ssl.use-sslv3	disable

so that the Nagle algorithm was disabled on all listen sockets to assure that there were no additional delays in total connection time. The default configuration parameters of the lighttpd server were used, except those listed in Table 2, which assured that the lighttpd server could efficiently handle the HTTP request rates considered in the experiments and that TLS was always preferred. The *OpenSSL* library v1.0.1e was used for the TLS protocol. The preferred cipher suite was ECDHE-RSA-AES128-GCM-SHA256, while Intel’s AES New Instructions (AES-NI) set for hardware accelerated AES operations was utilised [5]. The lighttpd server used the FreeBSD kernel SCTP implementation, and the Web traffic generator used the Linux kernel SCTP implementation.

An experiment run lasted for 600 s, during which the Web traffic generator generated exponentially distributed HTTP requests with a fixed average intensity, and with requested Web object sizes of 1 KiB and 35 KiB. In our experiment, we considered HTTP-request intensities ranging between 100 requests/s and 1000 requests/s. We measured:

- the total CPU load on the server,
- the CPU utilisation of every process that has a substantial contribution to the total CPU time,
- the total kernel memory used for networking.

Per-process CPU utilisation was sampled every 20 s and was calculated based on the accumulated CPU time given by

Table 3: Malloc types and zones

Command	Malloc type / Zone
vmstat -m	filedesc, kqueue, ip6opt, ip6ndp, pcb, BPF, ifnet, ifaddr, ether_multi, ltable, routetbl, igmp, in_mfilter, in_multi, ip_moptions, sctp_map, sctp_stri, sctp_stro, sctp_a_it, sctp_atcl, sctp_atky, sctp_athm, sctp_vrf, sctp_ifa, sctp_ifn, sctp_timw, sctp_iter, sctp_socko, hostcache, in6_mfilter, in6_multi, ip6_moptions, mld, inpcbpolicy, ipsecpolicy
vmstat -z	KNOTE, socket, udp_inpcb, udp_pcb, tcp_inpcb, tcpcb, tcptw, syncache, hostcache, sackhole, tcp_reass, sctp_ep, sctp_asoc, sctp_laddr, sctp_raddr, sctp_chunk, sctp_readq, sctp_stream_msg_out, sctp_asconf, sctp_asconf_ack, selfd
netstat -m	mbufs, mbuf clusters, 4k jumbo clusters, 9k jumbo clusters and 16k jumbo clusters

*procstat -r*. The total CPU load on the server was measured by measuring the accumulated CPU time of the *idle* system process (i.e., the total time that the CPU was idle) and subtracting this time from the total available CPU time during the measured interval (i.e., 160 s for an 8 parallel thread CPU). Total kernel memory utilisation was also sampled every 20 s and was calculated based on the output of *vmstat -z*, *vmstat -m*, and *netstat -m*. Table 3 outlines the malloc types and zones that were used to calculate total network-related kernel memory utilisation.

Our experiment comprised three test cases. In the first case, we evaluated a naive HE mechanism that did not employ caching of the outcome of previous happy eyeball invocations and which always resulted in a TCP connection being set up. The rationale behind this case was to serve as a baseline for the remaining two cases. Next, in the second test case, we still considered the same naive HE mechanism as in the first case, however, this time we evaluated happy eyeballing between TLS-encrypted TCP and SCTP. The second test case aimed at providing an appreciation of how the increase in CPU load and memory usage due to happy eyeballing compares with that caused by the TLS encryption itself. Lastly, in the third test case, we evaluated an optimized HE mechanism that employed caching of the outcome of previous connection attempts, using TCP and SCTP both with and without TLS encryption. The purpose behind the third test case was to obtain an understanding of the extent to which HE CPU load and memory usage decrease with caching, and to get a feel for the overhead of HE with a more optimized implementation. In this test case, we considered three different outcomes of the HE mechanism: HE always results in a TCP connection being set up (HE-TCP); HE always results in an SCTP connection being set up (HE-SCTP); and, HE results in a TCP connection being set up half the time and an SCTP connection half the time (HE-50%).

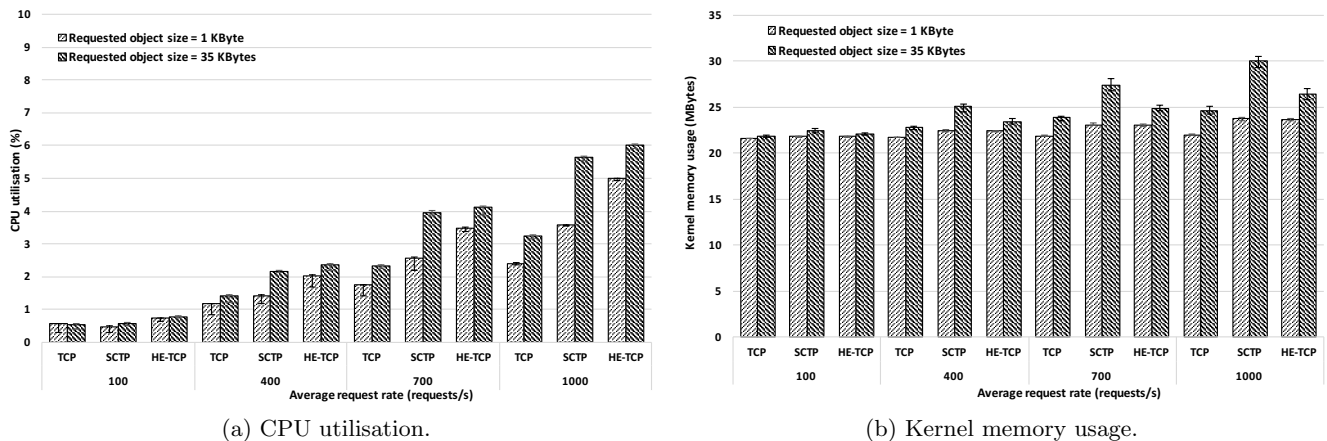


Figure 3: The results for the basic test case.

### 3.2 The Basic Test Case

The outcome of the basic test case is shown in Figure 3. Figure 3(a) compares the CPU load of HE between TCP and SCTP with that of single TCP and SCTP connection requests in the basic case. The figure shows how the CPU load varied as a function of the connection request rate and the size of the requested Web objects. The bar charts show the median values measured as described in Section 3.1, with error bars spanning the 10th and 90th percentiles. As follows, the CPU load of HE was quite substantial in the 1-KiB tests, roughly 40% higher CPU load than SCTP (i.e., the transport protocol considered in the study that consumed the most CPU load, in the tests with request rates 1000, 700, and 400 requests/s). However, as is evident from the 35-KiB tests, this was most likely an effect of the small amount of bytes transmitted in each Web response, and thus the small amount of bytes over which the CPU load was amortized: In the 35-KiB tests, the CPU load of HE was less than 10% higher than that of SCTP in the tests with request rates of 1000, 700, and 400 requests/s.

Figure 3(b) examines the kernel memory usage of HE compared with that of single TCP and SCTP connection requests. The bar charts show the median values with error bars spanning the 10th and 90th percentiles. Similar to Figure 3(a), the bar charts illustrate how the kernel memory varied with increasing connection request rates, and for different sizes of the requested Web object. We observe that HE had no or negligible impact on the kernel memory consumption – neither in the 1-KiB tests nor in the 35-KiB tests do we see a significant increase in kernel memory usage. In fact, the 35-KiB tests indicate that as the connection request rate increases, HE (at least in those cases where TCP wins) reduces the kernel memory usage as compared with SCTP.

### 3.3 Happy Eyeballing in the TLS Test Case

Figure 4 summarises the results from the TLS test case. Figure 4(a) is similar to Figure 3(a), but compares the CPU load in the case with TLS-encrypted connections. We observe that contrary to the basic case, the impact on CPU load of HE as compared with SCTP decreases significantly in the 1-KiB tests (less than 13% in all cases) and is not statistically significant in the 35-KiB tests (less than 4% in

the tests with request rates 1000, 700, and 400). Similar observations also apply when compared with TCP, where the impact of HE on CPU load is significantly lower than that in the basic case. Again, the reason HE had less effect on CPU load in this scenario compared with the basic case, was an effect of the way the CPU load was amortized.

Figure 5 illustrates how the CPU was shared among the kernel (including the FreeBSD subsystem) and the lighttpd server when the HE mechanism is used in the 35-KiB tests; both the tests for the basic case, as well as those for the TLS case. We observe that since the CPU load inflicted by HE was almost the same in both test cases, the CPU load of TLS (as reflected in the increase on the user CPU time of the lighttpd processes) overshadowed that of HE. Thus, in sum, we draw the conclusion that although HE is done at the price of some extra CPU load, the price becomes marginal for larger Web object sizes, and becomes even less significant in those cases HE is done between encrypted connections.

As regards the kernel memory usage in the TLS case, it follows from Figure 4(b) that HE had a marginal impact on this factor in this test case as well: In all tests, the kernel memory usage of HE is slightly higher than that of TCP (less than 8% higher memory usage), and always less than the kernel memory usage of SCTP.

### 3.4 Happy Eyeballing with Cached Results

In the basic and TLS use cases, we evaluated a naive HE mechanism that always tried both TCP and SCTP. This is, however, a rather inefficient implementation of HE. A more efficient and, as we see it, more realistic implementation would cache the outcome of previous connection attempts. So, e.g., assume that we have a cache hit rate of 80%, then HE tries both TCP and SCTP in only 20% of the application connection requests; in the remaining 80% of the application connection requests, HE issues either a TCP or SCTP connection request depending on the content of the HE cache. A cache hit rate of 80% is actually not an unreasonable figure. In statistics we obtained from Mozilla, they observed a hit rate of  $\approx 84\%$  in the Firefox internal ‘route’ cache during a six-week observation period.

Figure 6 shows the median CPU load (with error bars spanning the 10th and 90th percentiles) of HE at different

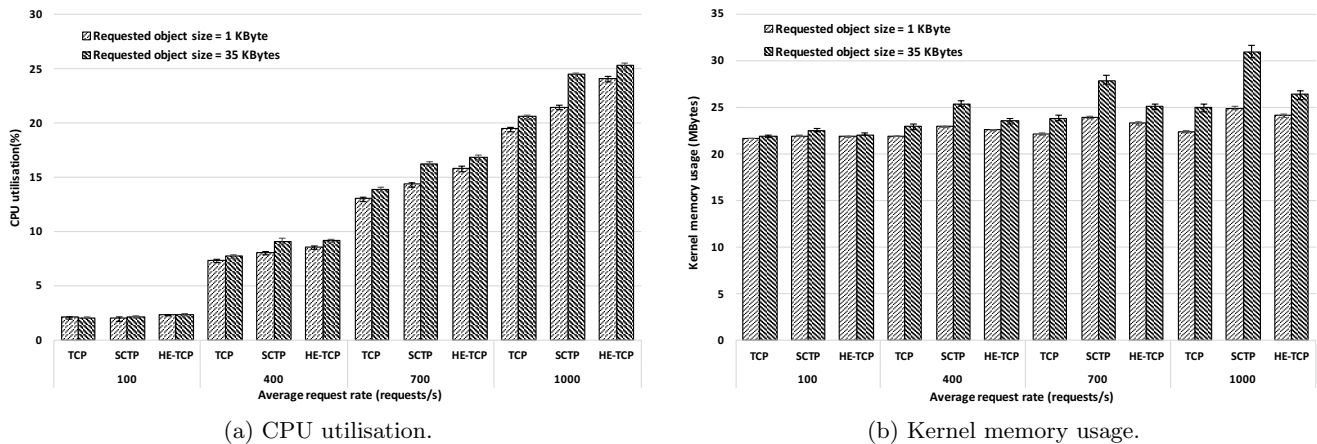


Figure 4: The results for the TLS test case.

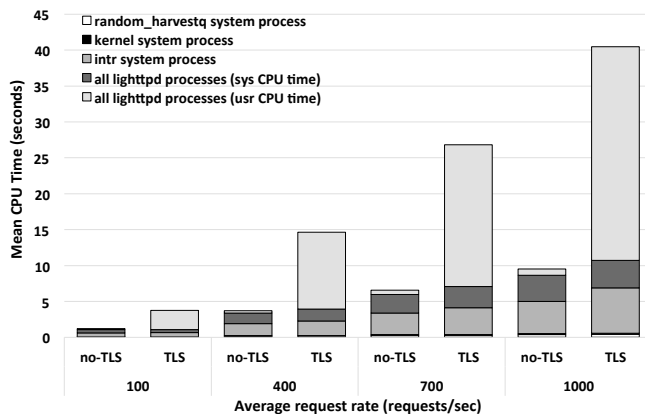


Figure 5: Breakdown of CPU utilisation for happy eyeballing between TCP and SCTP. Requested object size is 35 KiB.

cache hit ratios between 0% (naive HE that always tries both TCP and SCTP) and 100% (single TCP flows) in the 1-KiB tests: Figure 6(a) shows the total CPU load when TLS is not used, and Figure 6(b) when TLS is used. Figure 7 shows the corresponding results for the 35-KiB tests. Since the CPU load increases with increasing connection request rates and is thus more pronounced at higher request rates, we only consider the tests with a request rate of 1000 requests/s in Figures 6 and 7. Still, it should be noted that similar results were obtained for the lower connection request rates. The HE tests considered three outcomes: HE always results in a TCP connection being setup (HE-TCP); HE always results in an SCTP connection being setup (HE-SCTP); and, HE results in a TCP connection being setup half the time and an SCTP connection half the time (HE-50%).

We observe that the CPU load of HE decreases linearly as the cache hit rate increases, and this decrease is higher, percentage-wise, when TLS is not used (in the 1-KiB tests about 43% reduction in the CPU load when the cache hit rate is 80% and TLS is not used, and 18% reduction for TLS encrypted connections and the same cache hit rate). Again,

the reason caching had less effect on the decrease of the CPU load when TLS was used, was the effect of the way the CPU load was amortised. We further observe in Figure 6 that irrespective of whether TLS is being used or not, the difference in the CPU load imposed by HE-TCP, HE-50, and HE-SCTP is negligible for the 1-KiB tests. This implies that for small objects the additional cost of the http transaction is almost the same for both TCP and SCTP. For larger objects (e.g., 35 KiB), however, this cost is higher when SCTP is used, and hence significant differences between HE-TCP, HE-50, and HE-SCTP are observed at low cache hit rates in Figure 7.

We omit showing how the cache hit ratio influences the kernel memory consumption, since already the basic and TLS use cases suggest that kernel memory usage is not much of an issue for HE. Still, for completeness, we can mention that the cache hit ratio also had a positive impact on kernel memory usage. For instance, in the 1 KiB tests the kernel memory usage of HE at a cache hit rate of 80% was pretty much the same as for single TCP flows.

## 4. CONCLUSIONS

The Happy Eyeballs algorithm was originally proposed, and is currently being deployed, as a way of making a smooth transition from IPv4 to IPv6. However, the algorithm has also been proposed as a transport-selection mechanism. This paper evaluates happy eyeballing between TCP and SCTP, and shows that although HE increases CPU load as compared with a single TCP or SCTP connection establishment, the increase is in the order of 10% for 35 KiB Web objects, i.e., fairly typical Web objects, and is even smaller in those cases the happy eyeballing takes place between TLS-encrypted connections. Moreover, we show that the caching of connection-request results substantially reduces the HE CPU load, especially in comparison with the cost of TLS. As regards memory usage, our results suggest that HE has essentially the same memory footprint as single TCP/SCTP flows. The analysis in this paper shows that integrating a Happy Eyeballs mechanism into a library which provides a generic transport service is indeed a viable option to enable the use of advanced transport-protocol features whenever they are available. An example of such a library is the `neat`



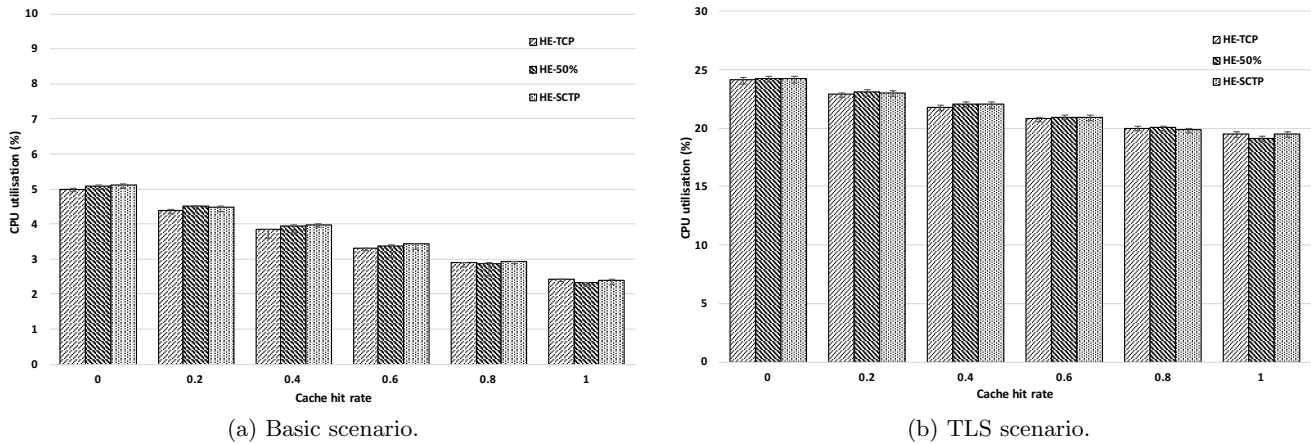


Figure 6: Impact of cache hit ratio on CPU utilisation. Requested object size is 1 KiB with a request rate of 1000 requests/s.

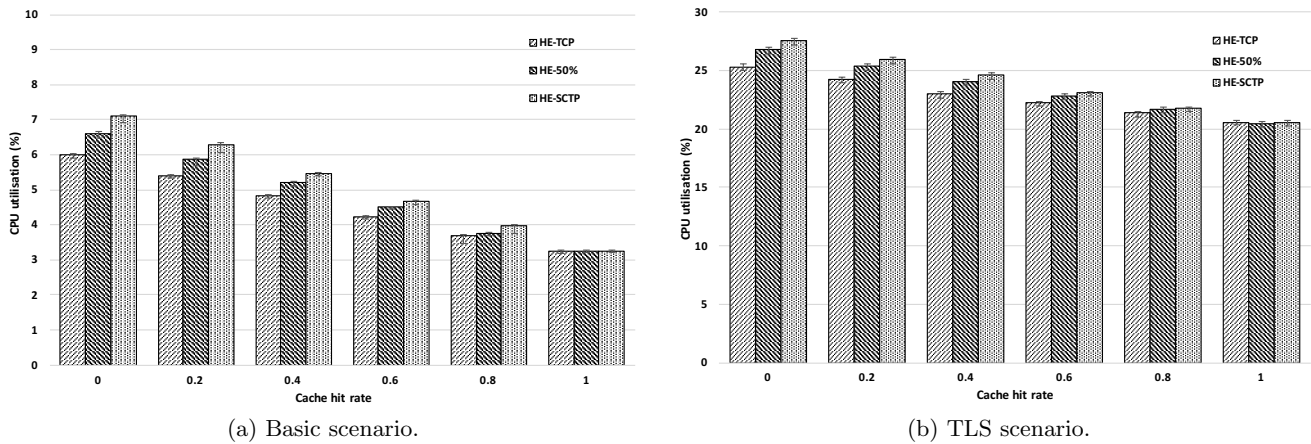


Figure 7: Impact of cache hit ratio on CPU utilisation. Requested object size is 35 KiB with a request rate of 1000 requests/s.

library currently being developed by the NEAT project<sup>2</sup>.

Future work includes evaluating the performance of HE when there are more than two transport solutions to be tried, e.g., TLS or DTLS encrypted traffic using IPv4 or IPv6 as the network layer, and TCP, native SCTP, or UDP-encapsulated SCTP as the transport layer, giving a total of six protocol candidates. Already, we note that the use of caching becomes even more important in these cases. Future work would also consider real-world experiments, where middlebox interference can be taken into account, as well as additional metrics to further examine the effects of transport happy eyeballing on both the network and destination end hosts, such as resource consumption on middleboxes, network load, and transaction times.

## 5. ACKNOWLEDGMENTS

The authors would like to thank Patrick McManus (Mozilla) for providing the Firefox cache-hit statistics.

<sup>2</sup><https://github.com/NEAT-project/neat>

This work has received funding from the European Union’s Horizon 2020 research and innovation programme under grant agreement No. 644334 (NEAT). The views expressed are solely those of the authors.

## 6. REFERENCES

- [1] E. Aben. Hampering Eyeballs – Observations on Two “Happy Eyeballs” Implementations. RIPE NCC, Nov. 2011. <https://labs.ripe.net/Members/emileaben/hampered-eyeballs>.
- [2] V. Bajpai and J. Schoenwaelder. Measuring the effects of happy eyeballs. Internet Draft draft-bajpai-happy, work in progress, July 2013. <https://tools.ietf.org/html/draft-bajpai-happy>.
- [3] F. Baker. Testing Eyeball Happiness. RFC 6556 (Informational), Apr. 2012.
- [4] O. Bonaventure. Happy eyeballs makes me unhappy..., Dec. 2013. <http://perso.uclouvain.be/olivier.bonaventure/blog/html/2013/12/03/happy.html>.

- [5] S. Gueron. Intel® Advanced Encryption Standard (AES) New Instructions Set. *Intel Corporation*, 2012.
- [6] httpperf. The httpperf page on SourceForge. <https://sourceforge.net/projects/httpperf>.
- [7] G. Huston. Bemused Eyeballs: Tailoring Dual Stack Applications for a CGN Environment. The ISP Column, May 2012. <http://www.potaroo.net/ispcol/2012-05/notquite.html>.
- [8] IETF. Transport Services (taps) Working Group. <https://datatracker.ietf.org/wg/taps/charter/>.
- [9] Lighttpd. Lighttpd – fly light. <https://www.lighttpd.net>.
- [10] D. Schinazi. Apple and IPv6 — Happy Eyeballs. Email to the IETF v6ops mailing list, July 2015. <https://www.ietf.org/mail-archive/web/v6ops/current/msg22455.html>.
- [11] D. Wing and A. Yourtchenko. Happy Eyeballs: Trending Towards Success (IPv6 and SCTP). Internet Draft draft-wing-tsvwg-happy-eyeballs-sctp-02, work in progress, Oct. 2010. <https://tools.ietf.org/html/draft-wing-tsvwg-happy-eyeballs-sctp-02>.
- [12] D. Wing and A. Yourtchenko. Improving User Experience with IPv6 and SCTP. *The Internet Protocol Journal*, 13(3), Sept. 2010. <http://www.cisco.com/c/en/us/about/press/internet-protocol-journal/back-issues/table-contents-49/133-he.html>.
- [13] D. Wing and A. Yourtchenko. Happy Eyeballs: Success with Dual-Stack Hosts. RFC 6555 (Proposed Standard), Apr. 2012.
- [14] D. Wing, A. Yourtchenko, and P. Natarajan. Happy eyeballs: Trending towards success (IPv6 and SCTP). Internet Draft draft-wing-http-new-tech, work in progress, Aug. 2010. <https://tools.ietf.org/html/draft-wing-http-new-tech-01>.

### **Disclaimer**

The views expressed in this document are solely those of the author(s). The European Commission is not responsible for any use that may be made of the information it contains.

All information in this document is provided “as is”, and no guarantee or warranty is given that the information is fit for any particular purpose. The user thereof uses the information at its sole risk and liability.