

Network Attachment Point (NAP)

Sebastian Robitzsch <sebastian.robitzsch@interdigital.com>

30th January 2017

Contents

1	Introduction	5
1.1	Installation and Configuration of IP Service Endpoints	5
1.1.1	Internet Connectivity through an ICN Gateway	6
1.1.2	Configure IP Service Endpoints	7
1.2	Configuration Variables	7
1.2.1	bufferCleanerInterval	7
1.2.2	fqdns	7
1.2.3	interface	8
1.2.4	icnGwNetworkAddress and icnGwNetmask	8
1.2.5	ipEndpoint	8
1.2.6	ltpInitialCredit	8
1.2.7	ltpRttListSize	8
1.2.8	ltpRttMultiplier	8
1.2.9	httpHandler	9
1.2.10	httpProxyPort	9
1.2.11	molyInterval	9
1.2.12	networkAddress and netmask	9
1.2.13	nodeId	10
1.2.14	routingPrefixes	10
1.2.15	socketType	10
1.2.16	surrogacy	10
1.2.17	tcpClientSocketBufferSize	11
1.2.18	tcpServerSocketBufferSize	11
1.3	Logging and Debugging	11
2	ICN Namespaces	13
2.1	Management	13
2.1.1	DNS Local	13
3	Transport Protocols	14
3.1	Unreliable Transport Protocol	14
3.2	Lightweight Transport Protocol	15
3.2.1	Header Format	15

3.2.2	Round Trip Time	16
4	Interfaces	17
4.1	NAP-SA	17
4.1.1	NAP_SA_ACTIVATE	17
4.1.2	NAP_SA_DEACTIVATE	18
5	Code Structure and Implementation Design Choices	19
5.1	Type Classes	19
5.1.1	Eui48	19
5.1.2	IcnId	19
5.1.3	IpAddress	20
5.1.4	Netmask	20
5.1.5	NodeId	20
5.1.6	RoutingPrefix	20
5.2	Demultiplexing Incoming Packets into their Handlers	20
5.3	Transparent HTTP Proxy	20
5.4	TCP Socket Handling	21
5.5	Network Attachment Point Packet Buffers	23
5.5.1	Buffer Cleaners	23
5.5.2	IP Packets Issued by Endpoints towards cNAPs	24
5.6	Co-incidental Multicast Group Formation	25
5.7	Traffic Control	26
5.7.1	Background	27
5.7.2	Implementation	27
5.7.3	Configuration	28

CID	Content Identifier
CMC	Co-incidental Multicast
cNAP	client-side Network Attachment Point
eNAP	extended Network Attachment Point
EUI-48	48-bit Extended Unique Identifier
FD	File Descriptor
FQDN	Full Qualified Domain Name
GW	Gateway
HTTP	Hypertext Transfer Protocol
ICN	Information-centric Networking
ID	Identifier
FID	Forwarding Identifier
LTP	Lightweight Transport Protocol
MAC	Medium Access Control
MITU	Maximum ICN Transmission Unit
MOLY	Monitoring Library
NACK	Negative Acknowledgement
NAP	Network Attachment Point
NID	Node Identifier
PID	Port Identifier
POINT	iP Over IcN - the betTer ip
rCID	reverse Content Identifier
RTT	Round Trip Time
RV	Rendezvous
SA	Surrogate Agent
SE	Session End
SED	Session Ended

SK	Session Key
sNAP	sever-side Network Attachment Point
STL	Standard Template Library
TCP	Transport Control Protocol
UE	User Equipment
URL	Uniform Resource Locator
UTP	Unreliable Transport Protocol
WE	Window End
WED	Window Ended
WU	Window Update
WUD	Window Updated

Chapter 1

Introduction

This document should be seen as some sort of Network Attachment Point (NAP) documentation to bridge the gap between the POINT deliverables [1], describing NAP functionality from a system perspective, and the code documentation available through Doxygen [5].

1.1 Installation and Configuration of IP Service Endpoints

The compilation and installation of the NAP has been tested on Debian-based Linux distributions, i.e.:

- Debian 8
- Ubuntu 14.x and 15.x
- Voyage 0.10

The following platforms have been successfully tested:

- x86 and x86_64 desktop and server machines
- APU and APU2s from PC Engines
- Raspberry Pi 3

The following virtualisation frameworks were successfully tested too:

- VirtualBox on Windows, Linux and MacOS hosts
- VMware on Windows hosts
- KVM on Debian and Ubuntu hosts

All required libraries in order to successfully compile the NAP are listed in the Blackadder-wide list of libraries. Simply install all of them before making the NAP:

```
~$ sudo apt install $(cat ~/blackadder/apt-get.txt)
```

As indicated in the README.md file, the NAP comes with a GNU-complaint make file which does not require any further customisation. Simply invoke `make` and provide the number of available cores to the `-j` argument to speed up the compilation, e.g. `make -j2` (no space between `-j` and the number of cores available for compilation]). To install the resulting binary called `nap` as a system-wide program run `sudo make install`. This copies the NAP binary to `/usr/bin` and creates the configuration file directory `/etc/nap` (if it does not exist yet) as well as a template configuration file directory, `/usr/share/doc/nap`.

When running the NAP binary, the required configuration files are all expected to be located in `/etc/nap`. When calling `sudo make install` the Makefile attempts to copy the template configuration files to `/etc/nap`. If the files already exist in the destination directory the user is prompted if they should be overwritten. Simply answer yes (`y`) or no (`n`).

1.1.1 Internet Connectivity through an ICN Gateway

When configuring the NAP as an ICN gateway towards the Internet the configuration file of the NAP acting as an ICN gateway must receive the following routing prefix configuration must be set (see Section 1.2.12 for more details):

```
networkAddress = "0.0.0.0";  
netmask = "0.0.0.0";
```

Furthermore, it is advisable to tell the ICN gateway the routing prefix which covers all routing prefixes configured in the ICN networks. This limits the number of packets the demux must process to the one which are targeted at IP endpoints attached to other NAPs (see Section 1.2.4).

The IP gateway which provides Internet access must receive the same configuration as IP service endpoints which do not have the NAP as their default gateway. Please see the next section how to configure them accordingly.

In addition to adding the ICN gateway routing prefix to the NAP which acts as the gateway all NAPs which are supposed to provide Internet access to their IP endpoints must receive the following routing prefix entry in their list of available prefixes (see Section 1.2.14):

```
networkAddress = "0.0.0.0";  
netmask = "0.0.0.0";
```

1.1.2 Configure IP Service Endpoints

In case the IP service endpoint attached to the NAP does not have the NAP as its default IP gateway (e.g. an IP gateway which is performing NAT or a web server which has two interfaces) the IP routing table must have an entry which basically says all traffic from IP endpoints attached to other NAPs must be sent to the NAP to which the IP service endpoint is attached to. Assuming the routing prefix which covers all NAPs is `172.16.0.0/16` and the NAP the IP service endpoint is attached to has the IP address `172.16.123.1` the following rule must be inserted into the IP routing table:

```
~$ route add -net 172.16.0.0 netmask 255.255.0.0 gw 172.16.123.1
```

Make sure that the interface which connects an IP service endpoint with its NAP has a subnet which does not cover any other IP endpoint attached to another NAP; this ensures that the configured gateway is always used for any communication and the IP service endpoint does not assume the IP endpoint is reachable within its subnet (this would cause ARP requests for the IP address which will remain unanswered).

1.2 Configuration Variables

This section explains the variables of the NAP's libconfig-based configuration file. The variables are listed in alphabetical order but can be placed in the configuration file in arbitrary order. All commented variables in the template `nap.cfg` file have the default value assigned to them so that the inexperienced user does not have to walk through the source code to find out which variable has been assigned with which default value.

1.2.1 `bufferCleanerInterval`

All handlers have a buffer in case a packet cannot be published under its Content Identifier (CID) due to various reasons (e.g. outstanding Blackadder notifications for this particular CID such as `START_PUBLISH` or `START_PUBLISH_iSUB`). The given value assigned to the variable determines the interval in seconds in which the buffer cleaner wakes up and checks all Information-centric Networking (ICN) buffer cleaners for packets older than the interval the cleaner wakes up.

1.2.2 `fqdns`

`fqdn` : String, mandatory

`ipAddress` : String, mandatory

`port` : Integer, optional

1.2.3 interface

This variable tells the NAP on which local network device it communicates with IP endpoints. The argument to this variable must be given as a string and the interface name provided must have the IP assigned all the NAP's IP endpoints send their IP traffic to. Reason being how the NAP utilises PCAP, i.e. reading the host IP address to compile a PCAP filter which ignores packets sent directly to the NAP.

Note, virtual interfaces and bridges have not been extensively tested and - so far - it can be stated that if the NAP is bound to a bridge this cause weird behaviour in the IP handler.

1.2.4 icnGwNetworkAddress and icnGwNetmask

More information about how to set up the NAP as an ICN Gateway (GW) is explained in Section 1.1.1.

1.2.5 ipEndpoint

For host-based deployments where the NAP servers a single IP endpoint only the following variable must be uncommented and the IP address of the IP endpoint the NAP servers is stated there. The value must be given as a string.

1.2.6 ltpInitialCredit

The lightweight transport protocol for Hypertext Transfer Protocol (HTTP) packet delivery is using a credit-based transport mechanism, similar to SPDY/HTTP2.

1.2.7 ltpRttListSize

For obtaining the Round Trip Time (RTT) for Lightweight Transport Protocol (LTP) timeout checkers the NAP keeps a list of measured RTT values to cope with potential RTT outlier values. This variable allows to configure the size of the list which is used to obtain the average of all reported RTT values.

1.2.8 ltpRttMultiplier

As explained in further detail in Section 3.2.2, whenever LTP starts a timeout counter to wait for a response from one of its receivers it uses a multiple of the previously measured RTT. This particular multiplier can be changed with the variable `ltpRttMultiplier` which accepts unsigned integer values.

1.2.9 `httpHandler`

In certain scenarios it is desired to not use the HTTP namespace for HTTP-level services. This can range from insufficient service level agreements to technical issues with particular HTTP services and the NAP being incapable to translate them properly into the namespace; or the content/service provider simply does not want to enable this enhancement. For those cases the HTTP handler can be turned off so that packets towards TCP Port 80 will not be mapped to the HTTP namespace anymore. Consequently, all traffic will be treated as pure IP and the IP-over-ICN namespace will be used. The respective variable `httpHandler` allows the boolean values `true` and `false` which turns the HTTP handler on and off, respectively.

1.2.10 `httpProxyPort`

For any HTTP traffic sent over TCP/IP with TCP destination Port 80 the NAP handles the traffic differently leveraging the HTTP-over-ICN namespace. As the NAP acts as a transparent proxy, a special iptables rule is inserted which forwards HTTP traffic to a port usually used by Squid¹, i.e. Port 3127.

1.2.11 `molyInterval`

The reporting of monitoring data points to the monitoring server is realised via Monitoring Library (MOLY), a library available through Blackadder. The interval of how often available data points are being sent off is up to the process (in this case the NAP). This interval can be configured here in seconds. The value given must be an integer. If the variable is not set or set to 0 the reporting to the monitoring agent is disabled. However, the NAP still calls the corresponding class to collect monitoring data from across NAP classes.

1.2.12 `networkAddress` and `netmask`

Each NAP acts in a particular routing prefix following the IP-over-ICN namespace definition. This routing prefix is configured here using a network address and a netmask. Both values must be given as a string.

If the NAP acts as an ICN GW the both variables must receive the following values:

```
networkAddress = "0.0.0.0";  
netmask = "0.0.0.0";
```

¹www.squid.org

Furthermore, it is advisable to manually configure the routing prefix the ICN GW is running on the interface towards the IP GW which provides access to the Internet. Please see Section 1.2.4.

1.2.13 `nodeId`

As long as the management API has not been finalised and pushed into a release branch of Blackadder, the NAP must know the node identifier of the ICN core it is running on in order to issue `publish_data_iSub` publications. The value must be given as an integer.

Note that this variable will soon be disappear, as Blackadder is going to provide management information (such as the Node Identifier (NID)) via a dedicated management API.

1.2.14 `routingPrefixes`

The routing prefixes available in within the ICN network can be configured using a list of pairs of `networkAddress` and `netmask`. As indicated in the example configuration file, both values must be provided as strings in a human readable format. The order of the prefixes does not matter, as the NAP will order them according to their size (as in how many hosts a particular prefix comprises).

1.2.15 `socketType`

First off, this option is not meant to be used unless there are issues with sent IP packets towards IP endpoints, i.e. they can be seen on the wire (with Tshark or TCPDUMP) but the endpoint does not reply or the NAP log states they have been sent off but nothing is seen on the wire. To date it seems that some Linux kernel/OS versions do not accept IP packets sent through a `IPPROTO_RAW` socket. To mitigate this problem, the NAP can switch to Libnet [3] as an alternative to raw Linux IP sockets. If `socketType` is not set or commented the NAP uses the raw IP socket implementation of Linux to send Ip packets to endpoints.

1.2.16 `surrogacy`

In case the NAP is supposed to allow the (de-)activation of surrogates, a listener must be started which allows a surrogate agent to activate/deactivate a particular surrogate server by utilising the NAP-SA interface. If this variable is uncommented and set to true the NAP opens a netlink socket and follows the NAP-SA interface specification accordingly.

1.2.17 tcpClientSocketBufferSize

When reading from a TCP a packet buffer must be created first which eventually determines the maximal TCP segment size to which the window size grows for larger file transfers. As the HTTP proxy distinguishes between TCP sessions towards clients and servers there are two separate variables which allow to configure the packet buffers.

1.2.18 tcpServerSocketBufferSize

See Section 1.2.17.

Note, for the time being the client-side Network Attachment Point (cNAP) does not perform flow control operations using LTP's Window Update (WU) and Window Updated (WUD) control messages. Hence, the Transport Control Protocol (TCP) server packet buffer must be equal or smaller than the LTP credit so that the cNAP never requires to enforce LTP flow control. The server socket buffer size can be easily calculated by

$$LTP_{credit} * MITU \leq BufSize \quad (1.1)$$

1.3 Logging and Debugging

The NAP utilises Apache's log4cxx logging library which allows a class-based logging combined with a highly customisable output formats. When installing the NAP a default log4cxx configuration file, `nap.l4j`, is placed in `/etc/nap` which has all available classes set to logging level `INFO`. The following logging levels are used in the NAP:

ERROR: A crucial error within the NAP which causes malfunction behaviour and/or a complete stop/crash/hang of a particular NAP functionality. If such an error occurs please double check your NAP configuration file or consult the POINT community on github.

WARN: An unexpected behaviour in the sequence of actions which causes the NAP to not being able to process the message properly.

INFO: Important information about the start/stop of a NAP module.

DEBUG: Debugging information about a particular class which informs the user about the overall functional healthiness of the NAP and the class where this logging level is set in particular.

TRACE: A per packet status while it traverses the NAP. Note, this can cause a sheer flood of logging messages. Please use with caution and enable only when debugging the NAP's internals.

The logging output can be configured to stdout and the filesystem using the `log4j.rootLogger` identifier. The default is set to both and the log file can be found in the default Linux log directory, i.e. `/var/log/nap.log`. The size and number of log files can be configured using the `log4j.appender.R` identifier. For more information on how to use log4cxx please consult Apache's documentation.

Chapter 2

ICN Namespaces

2.1 Management

2.1.1 DNS Local

The information item `DNSlocal` allows extended Network Attachment Points (eNAPs) to inform any other NAP about a change in the number of publisher for a particular Full Qualified Domain Name (FQDN) under the `/http` namespace. As all Forwarding Identifiers (FIDs) for `/http/fqdn` in cNAPs are only requested from RV/TM if they do not exist in the local ICN core, i.e. Blackadder, DNS local allows to trigger the flushing of FIDs in the core locally for a particular FQDN.

Chapter 3

Transport Protocols

3.1 Unreliable Transport Protocol

The Unreliable Transport Protocol (UTP) only serves a single purpose: fragmentation and stitching the packets back together at the other side. UTP is plugged to the IP handler and fragments all packets which exceeds the maximum acICN payload. If the NAP receives a packet larger than the Maximum ICN Transmission Unit (MITU) the packet gets fragmented by the sending NAP and put into an UTP packet as payload. Each UTP packet has a predefined header, as illustrated in Figure 3.1. The definition of the fields is defined by `utp_header_t` defined in `transport/unreliabletypedef.hh`.

Key	Payload Length	Sequence	State	Payload
-----	----------------	----------	-------	---------

Figure 3.1: UTP packet format

Key The key is generated by the sending NAP and allows to differentiate between several UTP sessions. The key is a sum of the hashed destination CID and a random value generated by `rand()`. This can be found in the `Unreliable::publish()` method in `transport/unreliable.cc`.

Payload Length This field indicates the length of the payload length followed after the UTP header.

Sequence This field indicates the position of the fragment for the packet reassembly after receiving all fragments.

State This field indicates if the packet is the first, an intermediate or the last one. In case the UTP packet carries a single fragment only without any

preceding or succeeding packet the state is set to "single" packet. The states are defined in the enumeration `TransportStates`.

Payload The payload of the UTP header carrying the IP packet.

3.2 Lightweight Transport Protocol

3.2.1 Header Format

3.2.1.1 Control Plane Header

The LTP packet header for control plane traffic is illustrated in Figure 3.2. There are seven different LTP control plane packets realised following the LTP specification described in Deliverable 2.3 of the POINT project[1]:

- LTP-Negative Acknowledgement (NACK)
- LTP-Session End (SE)
- LTP-Session Ended (SED)
- LTP-Window End (WE)
- LTP-Window Ended (WED)
- LTP-WU
- LTP-WUD

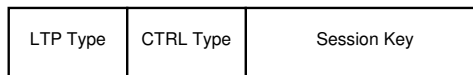


Figure 3.2: Lightweight Transport Protocol control plane packet format

All LTP control plane headers are implemented in `transport/lightweighttypedef.hh` as `struct` with LTP message type and control type messages enumerated in `enumerations.hh` (`ltp_message_types_t` and `ltp_ctrl_control_types_t`).

3.2.1.2 Data Plane Header

The LTP header for data plane traffic is illustrated in Figure 3.3 and consists of the four fields:

- LTP Type: Using `LTP_DATA` this field indicating that this LTP message is a data plane packet.

- Session Key: A per node unique integer identifying the HTTP session (derived from the socket file descriptor in the HTTP proxy)
- Sequence: A continuous sequence number indicating the position of the fragment
- Payload Length: The length of the LTP payload field in octets

LTP Type	Session Key	Sequence	Payload Length
----------	-------------	----------	----------------

Figure 3.3: Lightweight Transport Protocol data plane header format

3.2.2 Round Trip Time

RTT measurements are used as a timeout to discover that an LTP control message was potentially lost and must be therefore resent in order to keep the state machines in all NAPs participating in the same LTP session synchronised. The NAP measures RTT after finish publishing a full packet received from the IP endpoint by issuing an LTP WE packet and awaiting the corresponding response, i.e. WED. This behaviour is realised in the public `Lightweight::publish()` methods (`transport/lightweight.*`), one for HTTP requests with CID and reverse Content Identifier (rCID) and one for HTTP responses where the rCID and a list of NIDs is used. Both methods publish the data using `Lightweight::_publishData()` and issue a WU LTP control message right after. This is followed by a time-based counter to check for the corresponding awaited WED control message in order to proceed with the next packet from the IP endpoint. The time to wait is defined by a multiple of the currently known RTT. This multiplier is a fixed value stored in a private member of class `Lightweight`, i.e. `_timeout`.

To deal with measure RTT values much larger or smaller then the currently known RTT the NAP has a list of previously obtained RTTs and calculates the mean over them every time a new LTP session is created. This operation has been implemented in `Lightweight::_rtt()` which uses a default list size of 10 values. The list size can be configured using `ltpRttListSize` in the NAP's configuration file (see Section 1.2.7).

In order to accommodate for the possibility that a remote NAP disappears during an on-going LTP session the NAP always gives up to check for a received WUD after 23 attempts following the 23 enigma¹.

¹https://en.wikipedia.org/wiki/23_enigma

Chapter 4

Interfaces

This release of the NAP provides a single interface only, i.e. the NAP-Surrogate Agent (SA) interface.

4.1 NAP-SA

This section describes the communication interface between a SA and an NAP which allows to register a surrogate (authoritative delegate) at a NAP. The message exchange at this stage is the FQDN information that tells the NAP which FQDN the surrogate is serving and on IP address and port the IP service endpoint can be reached.

Both primitives to (de-)activate a surrogate server are listed next and are implemented in the class `NapSa` in `api/napsa.hh`. Based on the configuration variable `surrogacy`¹, the NAP opens a netlink socket with Port Identifier (PID) `PID_NAP_SA_LISTENER`, defined in `lib/blackadder_enums.hpp`.

4.1.1 NAP_SA_ACTIVATE

When generated This primitive is issued by SA whenever a static surrogate is fully operational and should be used.

Table 4.1: NAP_SA_ACTIVATION Primitive

Field	Type	Description
FQDN	<code>uint32_t</code>	The hashed version of the Full Qualified Domain Name
IP Address	<code>uint32_t</code>	The IP address in network byte order
Port	<code>uint16_t</code>	The port on which the IP service endpoint has opened its listening socket

¹See Section 1.2.16 for more information

Action upon arrival Using the scope path `/management/dnsLocal`, the surrogate NAP informs any other application subscribed to this particular scope path that the local FID entries must be flushed by unpublish and re-advertise the availability of information under `/http/hashedFqdn`.

4.1.2 NAP_SA_DEACTIVATE

When generated This primitive is issued by SA whenever a static surrogate is supposed to be removed from the network.

Table 4.2: NAP_SA_DEACTIVATION Primitive

Field	Type	Description
FQDN	uint32_t	The hashed version of the Full Qualified Domain Name
IP Address	uint32_t	The IP address in network byte order
Port	uint16_t	The port on which the IP service endpoint has opened its listening socket

Action upon arrival Using the scope path `/management/dnsLocal`, the surrogate NAP informs any other application subscribed to this particular scope path that the local FID entries must be flushed by unpublish and re-advertise the availability of information under `/http/hashedFqdn`.

Chapter 5

Code Structure and Implementation Design Choices

5.1 Type Classes

The NAP has implemented several helper classes to simplify the usage of certain types and the way they are used. Especially for logging purposes and ensuring the correct conversion between two representation of the same information the following type classes are used across the entire NAP code for EUI-48 addresses, ICN **IDs!**s (**IDs!**s), IP addresses, IP netmasks, NIDs and routing prefixes.

5.1.1 Eui48

The 48-bit Extended Unique Identifier (EUI-48) type is implemented in `types/eui48.hh` and allows the read and print of a 48 bit long address, commonly known as Medium Access Control (MAC) address. The conversation this class provides is simply between colon-free and colon-based string in- and output.

5.1.2 IcnId

The `IcnId` class which is implemented in `types/icnid.hh` is probably to richest type realised within the NAP. As the ICN Identifier (ID) is one of the key parameters when working with the Blackadder API the `IcnId` class simplifies the complexity reading, writing and changing the information of such an ID.

5.1.3 IPAddress

As the NAP's main purpose is the translation of packets coming from a standard IP stack into an ICN one the `IPAddress` class implemented in `types/ipaddress.hh` allows a hassle free read and write of an IP address (or network address) in network byte order or as a class-full representation.

5.1.4 Netmask

This class is very similar to the `IPAddress` class in the sense that it allows to read and write network bytes order and class-full representations. However, netmasks are also used in the CIDR notation where the netmask is provided as a single decimal number between 0 and 32. Thus, a dedicated class has been realised which can be found in `types/netmask.hh`.

5.1.5 NodeId

Blackadder (and the deployment tool in particular) has a strict way of accepting NIDs, i.e. an eight digit long string with leading zeros. In order to allow an easier way of following this convention the NAP has a `NodeId` class implemented in `types/nodeid.hh` which reads and writes both string and integer variants of the NID.

5.1.6 RoutingPrefix

The routing prefix utilises the implementation of the classes `IPAddress` and `Netmask` and wraps them up in methods to express a routing prefix in class-full, CIDR or network byte order formats. This class is implemented in `types/routingprefix.hh`.

5.2 Demultiplexing Incoming Packets into their Handlers

The demultiplexing of incoming packets slightly differs from the NAP's architectural description. While `iP Over IcN - the betTer ip (POINT) Deliverable 3.1` [2]

5.3 Transparent HTTP Proxy

As outlined in POINT deliverables [1], the NAP terminates TCP sessions aim to carry HTTP packets. Following IANA port assignments, all TCP traffic targeted at destination Port 80 is intercepted by the NAP and forwarded to the internal HTTP proxy port 3127¹. The proxy itself is imple-

¹see Section 1.2.10 on where and how to configure this port

mented in `proxies/http/httpproxy.hh` with two classes `TcpClient` and `TcpServer` to communicate with servers and clients, respectively.

Note, the two classes implement the functionality reflected in the class name and not for which type of TCP element they are serving. More precisely, the class `TcpServer` opens a listener socket on 3127 (if not configured differently) and once a TCP client tries to establish a TCP session with an IP service endpoint the NAP ()to which the client is attached to) intercepts the TCP session and the `TcpServer` class is called in a thread (see Line 161 and beyond in `proxies/httpproxy.cc`). On the contrary, an sever-side Network Attachment Point (sNAP) establishes a TCP session with a server if an HTTP packet has been received and needs to be send off to a particular IP service endpoint. For this the class `TcpClient` is used.

5.4 TCP Socket Handling

For HTTP-level services using TCP destination Port 80 the cNAP intercepts the TCP session and acts as a TCP server, as described in further detail in Section 5.3. Consequently, the sNAP acts as a TCP client towards the web server which serves the particular FQDN. As an HTTP session (i.e. request and its response) is very likely to be larger than a single TCP fragment and TCP socket reuse is used very often to reduce the number of opened sockets, it is important for the NAP to map the TCP socket state from an User Equipment (UE) attached to a cNAP to the sNAP which connects to the server for the duration of the entire HTTP session. This scenario is illustrated in Figure 5.1 which depicts UEs and the server in blue, and NAPs and their ICN links in red. Furthermore, the blue links between IP endpoints and the NAPs depict an exemplary socket file descriptor used by the respective NAP to communicate with a particular IP endpoint.

As explained in Section 3.2, LTP uses an Session Key (SK) to ensure the integrity of two HTTP sessions for te same web resource but requested by two UEs attached to the same cNAP and possibly with different HTTP headers (e.g., Range or User-Agent) which causes the web server to provide different HTTP responses. When accepting a new TCP connection from an UE at the cNAP the socket file descriptor becomes the SK which is part of LTP.

When the HTTP request is received by the sNAP (`PUBLISHED_DATA_iSUB` in `icn.*`) the LTP method `handle()` informs the callee if all segments have been received and therefore the received bytes (HTTP request) can be sent off to the server. At this stage LTP also returned the used SK so that the sNAP can hand this information together with the NID to the method `TcpClient::preparePacketToBeSent()` immediately followed by calling the functor `TcpClient::operator()()` to place the actual socket communication into a dedicated thread. Once this is done the respective

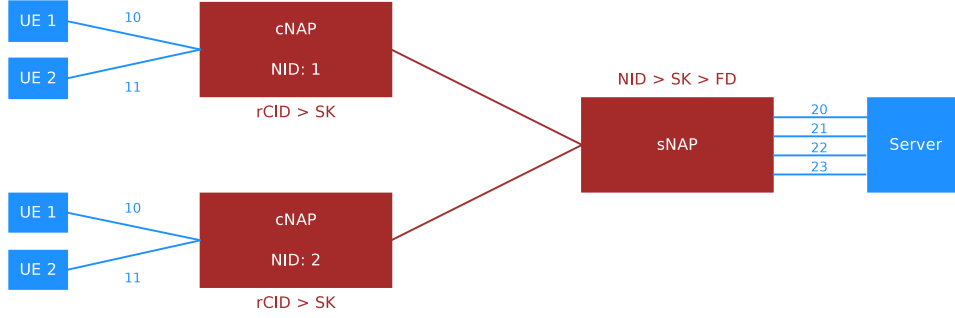


Figure 5.1: Handling of TCP sockets across Network Attachment Points an HTTP session

Table 5.1: TCP socket mappings across Network Attachment Points within private `_socketFds` map

NID	Remote Socket FD	Local Socket FD
1	10	20
1	11	21
2	10	22
2	11	23

thread looks up the private map `_socketFds` which holds a mapping of NIDs to remote socket File Descriptors (FDs) to local socket FDs, realised as an unordered Standard Template Library (STL) map within another STL map (`u_map<key, u_map<key, value>`). If a socket FD is found it means that a socket has been already opened and it can be re-used. In that way TCP socket re-use has been realised.

Table 5.1 depicts the mapping stored in `_socketFds` from the exemplary scenario illustrated in Figure 5.1. The NID is used as the key for the outer map and the remote socket FD as the key for the inner map (which is the value to the outer map's NID). Consequently, the local socket FD then is the value for the remote socket FD map key.

A mutex, `_socketFdsMutex`², is then used whenever an operation is performed on `_socketFds`, as the private members are shared among all threads created from the ICN handler class. So once a new HTTP request arrives at the sNAP the `_socketFds` map allows to look up if an existing socket FD is known; if not, a new socket is created. This functionality is implemented in `TcpClient::_tcpSocket()`. If the TCP client in the sNAP detects that the web server has shut down the TCP session or the socket is simply not readable anymore the local socket FD is getting removed from `_socketFds` map. If either the inner or the inner and the outer map are

²This Boost mutex is realised as a pointer to the respective class, as Boost does not allow to share the private mutex member being shared among all threads.

empty (no values left) they will be erased accordingly to keep the look-up time to find NID or remote socket FD keys to a bare minimum in the sNAP.

When the HTTP response issued by the web server is received at the sNAP it is potentially sent out via co-incidental multicast which means that all cNAPs which are in the Co-incidental Multicast (CMC) group will receive the response under the same randomly generated SK. That is why cNAPs keep the relation of published HTTP requests and their rCID and **PRID!** (**PRID!**) to the socket FDs which await the response. This is realised via the private member map `_ipEndpointSessions` in the class `HTTP` which upon arrival of an HTTP response at the cNAP the received rCID and **PRID!** is used to retrieve the list of UEs awaiting this response.

5.5 Network Attachment Point Packet Buffers

Another important implementation is the usage of packet buffers in order to accommodate for packet loss, enable packet reassembly and buffer packets to be published into the ICN network. In order to cover all the different buffers and implementation design choices to realise them this section is split into the following logical packet flows:

- Packets arrive from an IP endpoint at a cNAP
- Packets arrive at an sNAP from the local ICN core
- Packets arrive from IP endpoints at an sNAP
- Packet arrive at a cNAP from the local ICN core

The list above also reflects the logical flow of an IP communication (HTTP in particular). When handling IP packets in the NAP (using the IP handler) there is no difference whether the IP packet was issued by a server or by a client; in case this is an IP packet which carries HTTP it does make a difference.

5.5.1 Buffer Cleaners

In order to not buffer packets infinitely in case no subscriber ever appears the NAP has a dedicated buffer cleaner thread implemented which has access to both IP and HTTP buffers via pointers to the actual maps and their corresponding mutexes to guarantee thread safe read and write actions. The IP and the HTTP buffer cleaners are initialised in the IP and HTTP namespace class constructor, respectively. As the C++ namespacing of the NAP (to allow class-based logging³) does not allow to pass a reference to the IP and HTTP buffer and mutex the buffer constructors take void pointers only (see

³See Section 1.3 for more details about this functionality

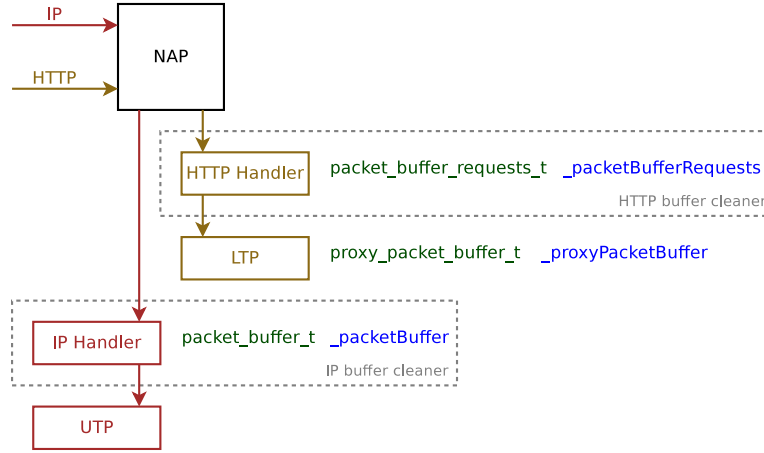


Figure 5.2: Packet buffers in the client-side Network Attachment Point for IP packets issued by endpoints

constructors in header files, `namespaces/buffercleaners/*.hh`) which are then casted back into their respective types (see constructor implementations in source files, `namespaces/buffercleaners/*.cc`).

A thread per ICN namespace is opened in each namespace constructor which wakes up every n seconds where n is the value of the NAP configuration variable `bufferCleanerInterval` (see Section 1.2.1).

5.5.2 IP Packets Issued by Endpoints towards client-side Network Attachment Points

This sub-section describes the used packet buffers which are used when IP packets were issued by an IP endpoint which behaves as an IP client in this scenario. A graphical representation of the used buffers can be found in Figure 5.2 with red boxes and arrows indicating IP traffic and golden boxes and arrows indicating HTTP packet handling.

IP For services other than HTTP (or any other future non-IP service) the IP handler is selected and the respective packet buffers to the left in Figure 5.2 are used. The implementation of the IP handler defined in `namespaces/ip.hh` uses a predefined type `packet_buffer_t` which uses the hashed CID as the map's unique key and a pair of the CID of type `ICnId` and a packet description struct of type `packet_t` as its values; both typedefs can be found in `namespaces/iptypedef.hh` and the `ICnId` class in `types/icnid.hh` (see Section 5.1.2 on Page 19). The buffer in the IP handler is solely used to buffer IP packets which could be published to the ICN core at the time. Reason being the entire CID needs to be

published first to the Rendezvous (RV) or the NAP has not received a **START_PUBLISH** notification for the CID under which the IP packet needs to be published. If a **START_PUBLISH** event arrives through the Blackadder API (see `icn.hh`) the ICN handler has a reference to all namespaces (i.e. `_namespaces`) which allows to look up the IP buffer for pending packets to be published via the method `Namespaces::publishFromBuffer` which switches based on the root scope into the particular namespace buffer (see `namespaces/namespace.cc`).

Any ready-to-be-published packet, either directly from the IP handler or through the buffer, is passed on to the UTP implementation which simply publishes the packet following the methods and procedures described in Section 3.1 on Page 14.

5.6 Co-incidental Multicast Group Formation

Whenever the NAP maps an incoming packet to the HTTP-over-ICN namespace, the formation of multicast groups is realised at the sNAP for the provisioning of HTTP responses to any cNAP which awaits such message. The logical steps required are illustrated in a message sequence chart in Figure 5.3.

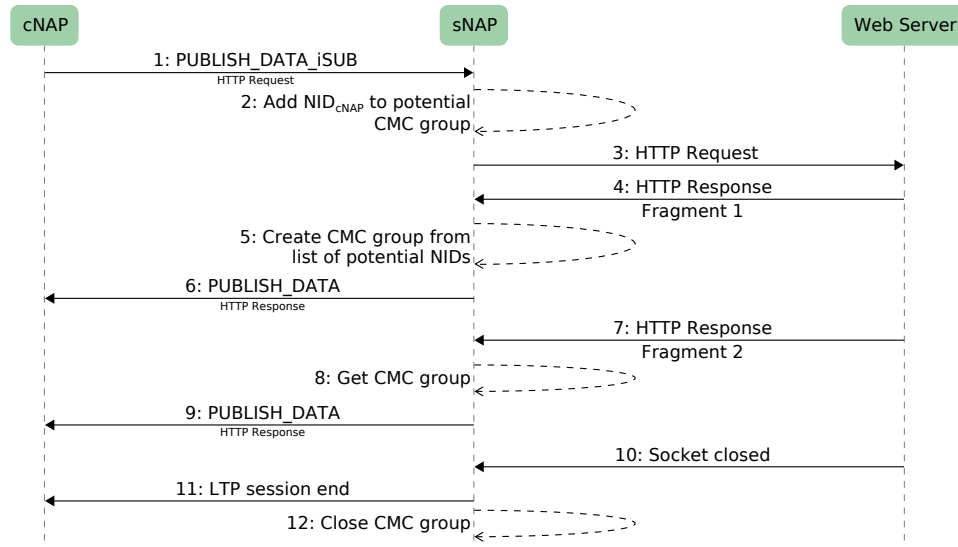


Figure 5.3: Co-incidental Multicast Group Management at sever-side Network Attachment Point

1. The cNAP has published an HTTP request to an sNAP which has subscribed to the respective FQDN.

2. Assuming there is no CMC group for this particular Uniform Resource Locator (URL), the sNAP creates a new potential CMC group with the hashed URL (rCID) as the unique key and adds the cNAP's NID to the list of NIDs in this potential CMC group.
3. The sNAP now delivers the HTTP request to the IP service endpoint via its transparent HTTP proxy (see Section 5.3 for more information about the proxy). Furthermore, it stores a backwards mapping of socket file descriptor to rCID.
4. The IP service endpoint responds with the first (and maybe only) fragment of the entire HTTP response. The HTTP proxy receives this via the socket created in Step 3.
5. The sNAP looks up the rCID for this particular socket file descriptor and obtains the rCID. With the rCID the sNAP finds all NIDs awaiting the response in the potential CMC group map. The sNAP closes the group from allowing future cNAPs to become a member.
6. The sNAP publishes the HTTP response to the CMC group of NIDs.
7. Assuming the HTTP response is made up of more than one fragment, the web server sends this to the sNAP which still has a TCP opened via its transparent HTTP proxy. Using the mapping from Step 3 the sNAP obtains the rCID and looks up if a CMC group is known which is awaiting this response.
8. The sNAP publishes the HTTP response to the group.
9. The web server closes the socket which is detected in the HTTP proxy
10. The sNAP closes the LTP session by sending an LTP-CTRL-SE message which must be confirmed by all **cNAPs!** (**cNAPs!**)
11. The sNAP closes the CMC group

5.7 Traffic Control

The NAP has received some preliminary functionality to perform traffic control actions on to-be-published data packets; it is important to understand that this functional extension is solely influencing packets originated from an IP endpoint. Blackadder-related control plane packets sent by the NAP, e.g., scope publications or publication of information items, are not affected by this.

5.7.1 Background

The reason behind this particular functionality was mainly caused by functional testing purposes of LTP (error control to be precise) in environments where the network does not cause any packet loss. Without any lost packet LTP's error control mechanism was hard to be tested and validated and using Linux's internal traffic control [4] would simply affect every single packet traversing a particular interface, including Blackadder control plane traffic. As no control plane reliability existed (at least at the stage when LTP was being tested), a `tc`-like module was written for the NAP which simply acts on packets to be publish using the `publish_data()` Blackadder primitive. As for `tc`, the respective traffic control module consists (theoretically) on the following filtering/traffic shaping mechanisms:

- Shaping
- Scheduling
- Policing
- Dropping*

However, at this stage only filtering/traffic shaping mechanisms marked with an (*) have been realised.

5.7.2 Implementation

The corresponding implementation of all `tc` mechanisms are collated in a derived class `TrafficControl` in `trafficcontrol/trafficcontrol.hh`. `TrafficControl` is then made a public member of class `Lightweight` in `transport/lightweight.hh`. Thus, all implemented traffic control gets only applied to packets published via LTP (for the time being).

In order to not break the code flow in the NAP when it comes to manipulating the publication of data packets the realisation of traffic control is realised as followed: every time a packet has been prepared to be published (across the entire LTP code in `transport/lightweight.cc`) the traffic control method `handle()` is called in an `if` statement. This method returns true if packet is supposed to be sent or false if not.

5.7.2.1 Dropping

The dropping of packets has been implemented as a binary decision within the class `Dropping`. When a packet is about to be published the respective LTP method calls `TrafficControl::handle()` from where `Dropping::dropPacket()` is called which has a boolean return.

The `dropPacket()` method first checks if dropping was requested by the user when adjusting/writing the NAP configuration file. This has been

realised through the method `Configuration::tcDropRate()` (declared in `configuration.hh`) which returns `-1` if `tcDroppingRate` has not been set in the configuration file. If that has not been the case C's standard element function `rand()`, from `stdlib.h`, is being used together with the configured drop rate:

```
rand() % _configuration.tcDropRate()
```

This essentially returns a value between 0 and `tcDropRate - 1`. If the outcome is 0 the method returns true (as in packet should be dropped); if the result is different from 0 the method returns false. The return value is directly returned to the place in LTP where `TrafficControl::handle()` has been called.

5.7.3 Configuration

5.7.3.1 Dropping

The dropping filter can be configured via the NAP configuration file using the variable

```
tcDroppingRate = <VALUE>
```

This variable accepts unsigned integer values between 0 and 2^{32} and represents the rate data packets will be dropped on average. For instance, if 100 is given to `tcDroppingRate` the NAP drops one in 100 to be published packets. Note, as LTP has been designed and implemented for well managed networks it is not advisable to set the rate lower than 100. Experiments have shown that in this case the LTP state machine of both publishing and subscribing endpoints gets out of sync.

Bibliography

- [1] POINT (iP Over IcN– the betTer IP).
- [2] Mays AL-Naday, Stavros Hadjitheophanous, George Petropoulos, Martin Reed, Janne Riihijärvi, Sebastian Robitzsch, Spiros Spirou, Dirk Trossen, and George Xylomenos. POINT D3.1 First Platform Design. Technical report, 2015.
- [3] George Foot and Peter Wang. Libnet.
- [4] Bert Hubert. Linux Advanced Routing and Traffic Control.
- [5] Dimitri van Heesch. Doxygen - Documentation Generator.