



Design and Implementation of a Deterministic Multipath Solution for Optimizing Backhaul Performance in Geographically Distributed 5G Campus Networks

Nicolas Zunker

n.zunker@campus.tu-berlin.de

January 3, 2024

MASTER'S THESIS

Telecommunication Systems Institute

Technische Universität Berlin

Examiner 1: Prof. Dr. Thomas Magedanz

Advisor: Dr. Ing. Marius Corici

Examiner 2: Prof. Dr. Axel Küpper

Declaration

I hereby declare that I have created the present work independently and by my own without illicit assistance and only utilizing the listed sources and tools.

Hiermit erkläre ich, dass ich die vorliegende Arbeit selbstständig und eigenhändig sowie ohne unerlaubte fremde Hilfe und ausschliesslich unter Verwendung der aufgeführten Quellen und Hilfsmittel angefertigt habe.

Die selbständige und eigenständige Anfertigung versichert an Eides statt:

Berlin, January 3, 2024

Nicolas Zunker

Acknowledgment

I would like to thank my unofficial advisor Hauke Buhr for his help, advice and support. I also thank my advisor Marius for his insights. Both of them provided excellent constructive feedback to me, for which I am grateful. I am also indebted to the rest team at Fraunhofer FOKUS for all their help, especially Hemant Zope, Eric Troudt, Christian Scheich and Aleksandar Yonchev. Lastly I would like to thank Prof. Dr.-Ing. habil. Falko Dressler and Prof. Dr. Zubow whose courses I enjoyed immensely during my studies, and from whom I learned a lot.

NZ

Abstract

The dawn of 5G brings with it an increased focus on bespoke quality of service provisions for different applications. Specifically for remote or emergency deployments, backhaul which respects the quality of service requirements of the applications running on-site (deterministic backhaul) is a difficult task. The new space race, in Lower Earth Orbit (LEO), has created a new possibility for backhaul over LEO satellite constellations and this transforms the backhaul challenge for industrial 5G sites, often located outside of cities and without access to fiber optic, into a multipath question.

This thesis prevents an investigation into, and proposal for, the use of multiple backhaul paths, in a geographically distributed 5G Campus deployment, to provide determinism. In the course of this document a design and an implementation is presented, which take a 5G flow's requirements and uses knowledge collected about the available backhaul paths to calculate, via a binary optimization equation, which path or paths to send it on. The proposed solution comprises both a control plane and a data plane component, as well as a traffic shaper, and a deployment configuration designed specifically for use in geographically distributed Campus 5G Networks is presented.

The described solution is investigated with respect to its ability to deterministically provide various network characteristics to traffic flows which require them, using a custom testbed which can emulate multiple links, . The results do not present conclusive evidence for the solution's ability to provide determinism, but do show promise. Although it is only able to achieve 80% optimality for latency based path switching, this thesis' solution is able to improve the latency experienced. It fails to reduce jitter for jitter-sensitive flows, and needs to compromise throughput in order to provide reduced packet loss, however it is able provide lower packet loss than even the Optimal Single Path Oracle, via its use of path switching and packet replication.

Zusammenfassung

5G legt einen größeren Fokus auf die verschiedene Servicequalität die von besonderen Applikationen gebraucht werden wird. Ins besonders abgelegene oder aus Notfall errichtete 5G Netze stehen vor einen schweren Problem, Backhaul bereit zu stellen der die benötigte Netzwerkeigenschaften ihrer kritischen Applikationen respektiert (deterministischer Backhaul). Dank der wachsenden Verfügbarkeit von Verbindungen über LEO Satellitenkonstellationen wird das traditionelle Problem vom Backhaul jetzt eine Multipath Frage.

Diese Arbeit beschreibt und implementiert eine Lösung die, in 5G Campus Umgebungen, verschiedene ausgehende Pfade benutzen kann um 5G Flows einen deterministischen Backhaul bereit zu stellen. Hierzu wird eine Architektur entworfen, welche eine binäre Optimierungsgleichung löst um diese Entscheidung zu treffen, anhand von Daten die über die ausgehenden Pfaden gesammelt werden. Des weiteren verfügt die Lösung über geteilte Steuerungs und Datenebenen und einen Traffic Shaper, und es wird eine Konfiguration spezifisch für 5G Campus Netze vorgestellt.

Die vorgestellte Lösung, der "WAN Connector", wird dann in einem Testbed der entworfen worden ist um Multipath Szenarien zu Emulieren auf ihrer deterministischen Eigenschaften untersucht. Die Ergebnisse weisen leider nach dass die Lösung nicht kompletten Determinismus anbieten kann. Obwohl es in der Frage der Latenz nur 80% Optimal agieren kann, schafft es, die Implementation dieser Arbeit, die Latenz insgesamt zu verringern. Den Jitter zu reduzieren gelingt nicht, und auch der Durchsatz von Flows die der WAN Connector für Paketverlust und Zuverlässigkeit optimiert wird extrem reduziert, dennoch ist es aber möglich die Paketverluste mehr zu verringern als sogar der Einzel-Pfad Orakel, durch die Nutzung von Pfad-wechsel und Replikation.

Table of Contents

Declaration	I
Acknowledgment	III
Abstract	V
Zusammenfassung	VII
List of Figures	XII
List of Listings	XIII
1 Introduction	1
1.1 Motivation and Problem Description	1
1.2 Goal	2
1.3 Structure	3
2 Background	5
2.1 Mobile Networks and Backhaul	5
2.1.1 Backhaul in 5G	8

2.1.2	5G and Campus Networks	9
2.2	Multipathing and Multihoming	10
2.2.1	Collecting Link and/or Path Performance Metrics	11
2.2.2	Path and/or Interface Selection	12
2.3	Determinism in Computer Networking	13
2.3.1	Deterministic Networking Specification	13
2.3.2	Traffic Shaping, Fair Queuing, and Bufferbloat	14
3	Approach	17
3.1	Requirements	17
3.2	Skeleton Structure	18
3.3	Architectural Components Required for Deterministic Multipath Backhaul in a 5G Campus Environment	19
3.3.1	Elimination of Contention Loss and Bufferbloat	20
3.3.2	Jitter Reduction and Latency Guarantees	22
3.3.3	Service Protection	22
3.3.4	Multipath Considerations	25
3.3.5	Summary of Necessary Features	26
3.4	Overview of the WAN Connector’s Features and Components	26
3.4.1	Path Selection Algorithm	26

3.4.2	Packet Ordering and Elimination Function	28
3.4.3	Internal Architecture	29
3.4.4	Overview	32
4	Evaluation	35
4.1	Approach to Evaluating Performance	35
4.1.1	Accuracy of the Emulation	38
4.2	Latency Based Path Switching	38
4.2.1	Performance Compared to Single Links	40
4.2.2	Latency Performance Under Load and the Importance of Traffic Shaping	42
4.3	Reliability Based Path Switching	44
4.3.1	Impact on Throughput	46
4.4	Jitter Based Path Switching	47
4.5	Discussion	49
4.6	Improvements	50
5	Conclusion and Outlook	55
5.1	Summary	55
5.2	Implications and Further Areas of Research	56
A	Patch for the BPF Flow Dissector to Peek Inside GTP	57

B	JSON Formats Used in Control and Data-Plane Messages	61
C	Example of Custom GTP Header Format	67
	Bibliography	69

List of Figures

1.1	5G Deployment with 2 UPFs	2
2.1	Simplified 5G Architecture	6
3.1	Internal Architecture of the WAN Connector	32
4.1	Testbed Setup	36
4.2	Comparison of the Emulators Latencies and the Testbed	39
4.3	Cumulative Distribution of the Latencies	41
4.4	Closeup of the Latency CDF	42
4.5	Latency with Background Traffic	43
4.6	Reliability	45
4.7	Throughput of a Duplicated Flow	46
4.8	Performance of Jitter Based Path Switching	47
C.1	Screenshot of Wireshark Dissection	68

List of Listings

A.1 Patch to remove GTP Headers for the BPF Flow Dissector Example in the Linux Kernel	58
B.1 Format of Flow Forwarding Decisions from Control Plane	63
B.2 Format of Statistics from Data Plane	65

Chapter 1

Introduction

1.1 Motivation and Problem Description

One of the aims for the fifth generation of mobile networks (5G) and its successors will be a greater diversification of the classes of service. As the use cases for these networks evolve, there is a greater need for quality of service (QoS) tailored to each use case. For example, in the Industrial Internet of Things (IIoT) the requirements on latency, jitter, and reliability may be extremely stringent. Supporting these kinds of classes of service can be a challenge for mobile network operators (MNOs) and will require novel approaches to familiar problems, such as backhaul.

As there are more heterogeneous edge deployments and more campus networks, backhaul becomes more challenging, since many sites may not have access to optical fibre, and may be forced into using other solutions such as satellite links, mmWave backhaul, or pre-existing on-site ISP connections. Providing the kind of deterministic quality of service that these sites may require can be a very difficult challenge.

Particularly with the rise of satellite backhaul options, fuelled by the new space race, many

remote deployments may choose to integrate satellite backhaul because fibre is infeasible or its rollout is too slow. These connections are usually being added in addition to existing one, and thus network operators may choose to utilize more than one backhaul connection at the same time. Either in order to increase the available bandwidth or to utilize the different qualities of the backhaul links. This bears the question whether multipathing could then be used to provide deterministic backhaul by intelligently selecting on which links to forward packets. This approach bears similarity to multihoming as well as to multi-path routing in Wireless Sensor Networks (WSNs), and can take inspiration from the substantial body of research in these fields which already exists, and which has demonstrated that QoS can be improved by using multiple links or paths simultaneously [6, 7, 17, 19, 21, 36].

1.2 Goal

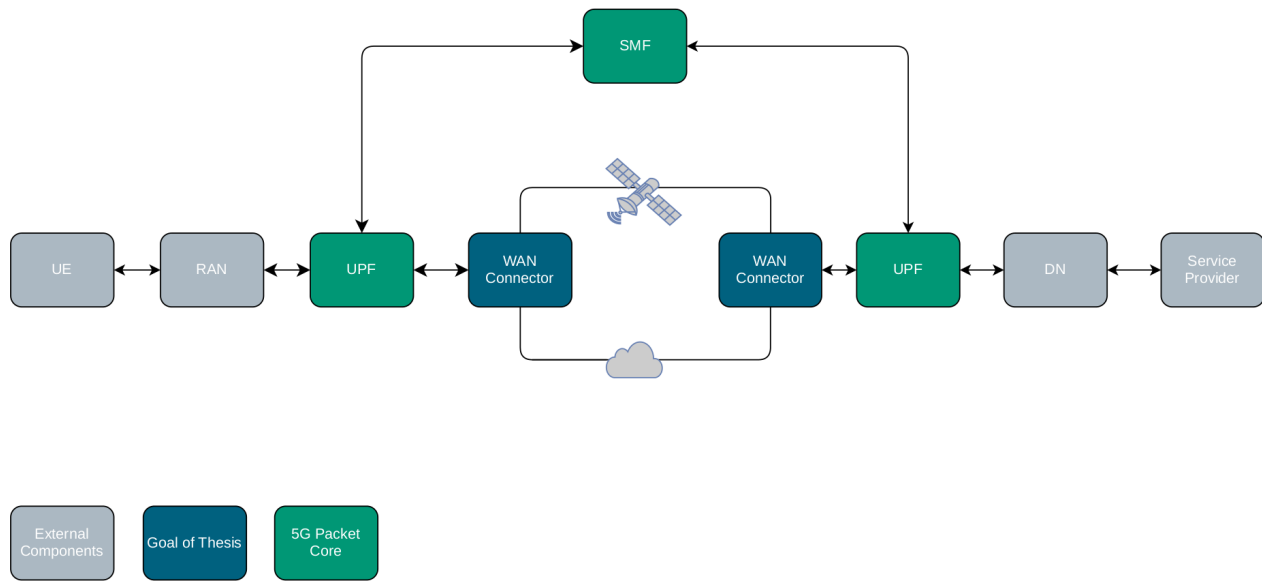


Figure 1.1: 5G Deployment with 2 UPFs

The goal of this thesis is to design and provide an implementation of a Wide Area Network (WAN) connector, that can be placed at the ingress and egress point of two or more locations, and utilizes multipathing in order to provide deterministic backhaul between the two sites.

The performance of this approach will then be quantitatively analyzed in experiments.

Looking at Figure 1.1 we can see how this is envisioned to work. WAN connectors are deployed both in the geographically distributed 5G campus network, which has more than one egress link, and in the core. Then, using the multiple outgoing links, the traffic is backhauled to the other site, while respecting its QoS requirements. This can be especially beneficial for critical applications (e.g. at industrial sites), which are in locations that do not have access to optical fibre for backhaul.

For a 5G deployment the proposed WAN connector could be deployed in between two User Plane Functions (UPF) , in order to provide deterministic backhaul. The architecture for such a deployment is shown in Figure 1.1. Further, the on-site UPF is not a strict requirement; it is also feasible to connect the RAN directly to the WAN Connector.

1.3 Structure

This thesis will follow a simple 5 chapter structure. This section concludes the introduction chapter, what follows will be one chapter to provide both basic background information as well as to highlight the existing literature which is of relevance to the problem statement. Next will be the approach and evaluation chapters, where, respectively, the design of the implementation will be presented and discussed, and then its performance will be analyzed. The conclusion chapter will review the relevant findings, the successes and failures of the approach, and the possibilities for future research in this area.

Chapter 2

Background

To place this thesis in its proper academic and practical context, this chapter will provide background information on relevant topics as well as reviewing some of the existing approaches in literature which bear similarity or served as an inspiration. First mobile networking and 5G will be discussed, then multipath solutions will be covered, and lastly deterministic networking.

2.1 Mobile Networks and Backhaul

For the sake of brevity, the full history of the development of mobile networks will be hugely condensed and simplified, and only the 5G architecture will be explained in depth.

The defining characteristic of a mobile network is wireless connectivity. This enables the users to move around while remaining connected to the network, which stands in direct opposition to historical computer networks which were stationary and usually featured fixed connections between computers. A user in a mobile network must first detect and then connect to a Base Station, this process includes authentication and authorization. A user may be switched

from one Base Station to another depending on the signal strength (this is referred to as a handover), and this is what enables the user to maintain a connection whilst moving across wide areas. In mobile networks which support Internet traffic, the user will normally appear remain behind a single IP address despite their mobility, and often behind a layer of Network Address Translation (NAT) as well.

All of the user's data sent in the network goes through the Base Station to which it is currently connected. At first this data was exclusively analog voice data, however over time newer versions of the mobile networking specifications have defined different types of data, and eventually defined support for Internet Protocol (IP) messages as well. Starting with 4G, mobile networks transitioned to an all IP architecture [1]. To support and manage user mobility and provide internet connectivity, today's mobile networks have a complex core network architecture where different network functions (NF's) are responsible for the different aspects of user management. This core network continues to evolve with each new specification.

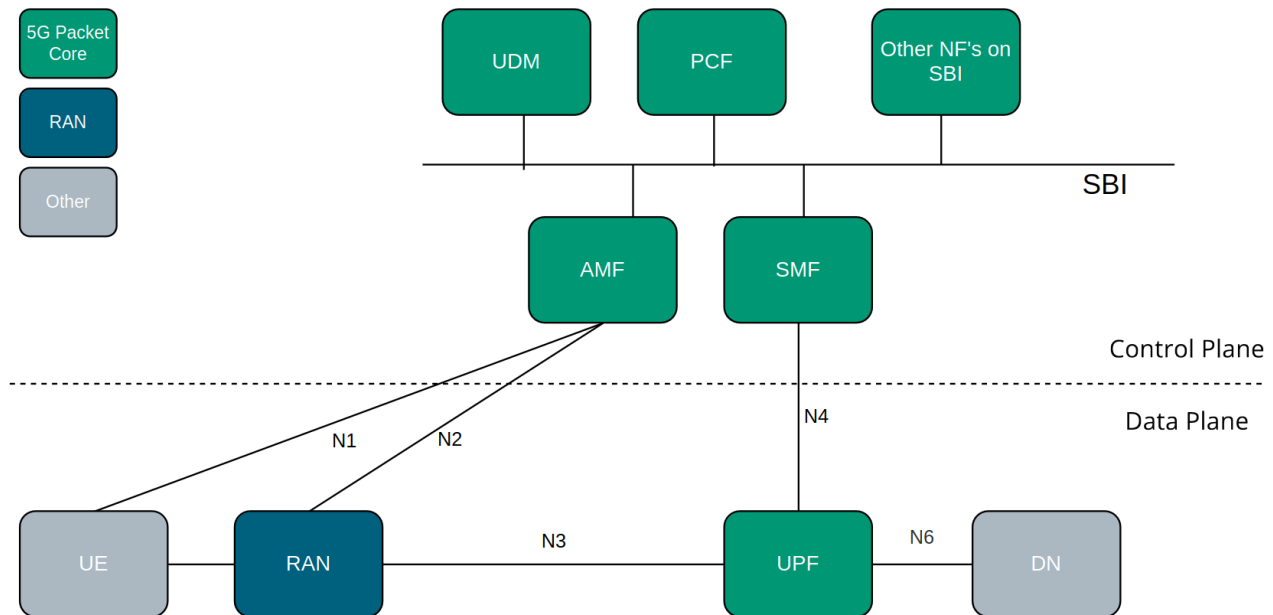


Figure 2.1: Simplified 5G Architecture

Figure 2.1 shows a simplified version of the architecture of a 5G core network and how it

is connects to the RAN, the user, and the internet. The figure does not include all of the Network Functions defined in the 5G specifications but shows a select few which will be briefly explained in following.

The Access and Mobility Management Function (**AMF**) handles, among other things, the user registration (access authentication and authorization), connectivity management, and the control plane interface to the RAN. The reference point for interactions between the AMF and the RAN is called the N2 reference point. The protocol used for these interactions is the NGAP protocol. The reference point for interactions with the UE is N1. This is a purely logical interface, messages between the UE and the AMF actually have to physically go through the RAN.

The User Plane Function (**UPF**) is responsible for forwarding packets between the Data Network and the RAN. The UPF must also perform QoS enforcement, for example making sure users do not exceed the limits of their subscription. The UPF communicates with the RAN over the N3 reference point. For this they utilize the GPRS Tunneling Protocol (GTP). GTP is a tunneling protocol designed for mobile networks. Packets originating in the RAN are contained in a GTP tunnel and must be de-capsulated by the UPF before being forwarded to the Data Network. Packets in the opposite direction must have a GTP header added to them before they are sent to the RAN.

The Session Management Function (**SMF**) is responsible for session establishment and management, and it configures traffic steering in the UPF. The SMF and the UPF communicate over the Packet Forwarding Control Protocol (PFCP) on the N4 reference point. Their interactions are a lot like those between a switch and a controller in a Software Defined Network, with the SMF acting as the controller.

The Policy Control Function (**PCF**) provides policy rules to different Control Plane functions in order to enforce them. The SMF consults the PCF before informing the UPF what Quality of Service a given user should experience. The interactions between the SMF and the PCF go

over the Service Based Interface (SBI). This is a logical interface which represents all of the different interactions between all of the Control Plane Network Functions. Communication on the SBI is done with the http2 protocol.

In the figure it can also be seen how the 5G Core (like the 4G Packet Core before it) splits the control plane and the user plane (or data plane). The network functions in the control plane co-ordinate to make decisions over how to orchestrate the data plane, while the components in the user plane are the ones forwarding the actual user traffic.

2.1.1 Backhaul in 5G

The first hop from the user equipment (UE) into the mobile network occurs over the air, between the UE and the Base Station, via wireless signals. The network of different Base Stations is called the Radio Access Network (RAN). Once the user has transmitted its data to the RAN, the RAN is responsible for passing on the user's messages to the core network. The transport network between the RAN and the Core is the backhaul network or backbone [22], however the term backhaul can also be used to refer to the process of transporting data from the RAN back to the core.

In light of the increased focus in 5G on enabling time sensitive industrial use cases, the authors in [28] proposed a method for integrating asynchronous Time Sensitive Networking into 5G backhauling. To achieve this their architecture uses an Asynchronous Traffic Shaper, which is crucial for integrating the time sensitive traffic with the regular traffic. Their solution also utilized frame replication for the delay critical flows defined in the 5G specifications (5G QoS Identifiers (5QIs) 82 through 85).

The authors of [22] highlight the need for innovation in 5G backhaul provisioning. At present optic fibre backhaul is not available nationwide in Europe, and Fibre to the Home (FTTH) is scarce worldwide. Currently most backhaul networks are built with microwave links or fibre/copper based links. There has also been research [32, 33] into using multipath

routing within the context of millimeter Wave (mmWave) backhaul in order to increase both reliability and throughput.

Another option for mobile backhaul is satellite constellation networks. Geostationary Orbit satellites (GEO satellites) appear motionless to observers from the ground, and thus provide very consistent latency and jitter. However this comes at a cost, they are a large distance away from the earth and round trip times are routinely on the order of 400 to 500 milliseconds. Regardless, GEO satellites are still a crucial means for connectivity for remote regions and areas where it is not possible or very difficult to provide other means of backhaul. Recently, Low-Earth Orbit (LEO) satellite constellations have become a more appealing option due to their greatly reduced latencies when compared with GEO satellites [12]. The 3G Public Partnership Project (3GPP), who are behind the 5G specifications, have commissioned a technical report, TR 23.737 [3], on possible architectures and integration scenarios for satellite backhaul in 5G.

2.1.2 5G and Campus Networks

In order to enable new industrial use cases, 5G campus networks may be deployed exclusively for users within a campus organization, and covering only a prescribed geographical area [31]. The deployment of such a network can be crucial for time-sensitive wireless communication, e.g. between robots, and the 5G specifications contain traffic classes and QoS requirements for delay-critical use cases [4], this is done in order to make 5G a competitive, or even superior, option to WiFi in campus deployments [39].

There is lots of room for different deployment architectures for 5G campus and non-public networks (NPN's). In [29] many of these architectures are discussed and investigated. In order to operate a 5G campus network, spectrum may need to be acquired, as well as 4G or 5G base stations, depending on the mode of operation- NSA (non-standalone: 4G Base Stations (eNodeB's) but a 5G Core Network) or SA (standalone: 5G Base Stations (gNB's)

and a 5G Core Network). Furthermore, a 5G Core Network is required in order to manage the network. 5G places an increased focus on virtualization and flexibility of deployment. As such there are a litany of options for which network functions from the core to deploy, and where to deploy them.

This thesis addresses specifically those campus network deployments which require backhaul to a core located off premises. Networks deployed in this fashion will typically feature a subset of the core's network functions deployed on premise. In the simplest case this subset will be the empty set, and the only component on premise will be the RAN. A more common deployment model may feature one UPF on site with a second UPF located in the core, and packets being forwarded between these two UPFs. The WAN Connector will be designed to be able to backhaul IP traffic from the remote campus to the core. This means it can be placed "after" the RAN or "after" the UPF in order provide deterministic backhaul for the IP traffic it is forwarding. Here "after" refers to the order in which network components are traversed, from the perspective of the uplink traffic. In the core, the terminating WAN Connector will be placed "before" the UPF.

2.2 Multipathing and Multihoming

At this point it becomes important to establish the definitions of many similar terms, and how they will be used in this thesis. Multihoming refers to a host which is connected to more than one Internet Service Provider (ISP). As such it has more than one outgoing interface on which it may forward packets for which it does not have a direct link to the destination. A multihomed host may choose to use different links for different traffic. Multipath routing, which may also be called "multipathing", is a routing approach which selects separate paths when routing packets. In order to perform multipath routing across the internet a host must be multihomed. If the destination host is also multihomed then there are a large number of possible paths to choose from.

This discussion directly ties into the difference between links (or interfaces) and paths. Especially in the example with two multihomed hosts, using the same outgoing link does not mean one is always the same path. The same interface could be used to send packets with different destination IP addresses (for the same host), but the packets would take different paths. In the context of a multihomed host communicating with a destination that only has one IP address then selecting a different interface is the same as selecting a different path. In such cases multihoming and multipathing are the same.

2.2.1 Collecting Link and/or Path Performance Metrics

In [7] the authors collected both passive metrics (looking at response times for outgoing packets), and active measurements (sending ICMP ping, or TCP SYN messages and measuring the response time). Using the passive measurements enabled their multihomed approach to perform well, but when using the active measurements the performance was even better. Crucially, the passive measurements worked better over larger sampling periods, because it took longer to get a full overview of all the possible routes. Whereas the active sampling approach acquired its measurements faster and was thus more effective over smaller sampling intervals.

Considering this result, the solution developed in this work will therefore utilize both active and passive measurements. All three metrics- packet loss, latency, and jitter- will be periodically measured in an active manner. The period over which to perform these measurements is also an important design decision for the WAN connector, and it will be a fixed value, which can be configured by the operator. The use of a fixed value is important in the context of a solution which aims to provide determinism, since traffic may often have hard requirements for the maximum tolerable latency, packet loss, or jitter, which may be experienced within a given window (e.g. in the 5G QoS specifications the averaging window for non delay critical traffic is 2 seconds [4]).

For classical wired links in multihomed scenarios, [35] have observed that one link will generally dominate with regards to latency, but with brief periods where other links' performance is superior. These same authors also note that with regards to packet loss there is far less prevalence of a "dominant" link, and generally the links will perform comparably. Packet loss is also a particularly difficult metric to measure, since most links are highly reliable and when they do experience packet loss it is in bursts [35]. Wireless connections are usually less reliable and may experience more consistent rates of packet loss at the data link layer, however often this is opaque from the perspective of the higher layers and they may only perceive it through the jitter and/or latency.

Ultimately, many of the best performing approaches for predicting packet loss, e.g. Hidden Markov Models [11, 35], are still somewhat imprecise and inaccurate. These models assume the link is in one of two states, good or bad, and each state has a different probability for packet loss, and there is a transition matrix which represents the probabilities of switching from one state to another.

2.2.2 Path and/or Interface Selection

Multipath routing approaches, as well as interface selection in multi-homed environments are mature and well researched topics. Here, a selection of such approaches will be covered.

Wireless Sensor Networks (WSNs) are often deployed in challenging environments and must construct their routing tables dynamically. Research into multipath routing solutions for providing QoS guarantees in WSNs can serve as inspiration for this thesis' approach. In [9], the authors use Routing Request (RREQ) flooding to build knowledge of disjoint paths. The RREQ's also include a minimum reliability so that the paths found all comply to this requirement. Then when maintaining routing table entries the probability of success is also stored. This information is then used in the path selection. While this approach is good for WSNs, in the context of a 5G campus network it is not feasible because it uses the Ad-hoc

On-demand Distance Vector (AODV) Routing [27] protocol, whereas the public internet uses BGP [30].

A more relevant example is the work in [21], where the authors also address QoS provisioning in a WSN, using multipath routing. Their approach is to formulate the problem as a probabilistic programming problem and then convert it to a linear programming problem using an approximation technique. When solved, the equation they use to represent the problem yields which paths packets should be sent on (or duplicated on) in order to achieve the desired reliability. Only those paths which would guarantee the maximum latency is still below the QoS requirements are considered. This approach bears similarity to how [28] (discussed earlier) replicates certain flows.

Another article which deserves mentioning is the work in [7]. Here, the authors evaluate a large set of route-control mechanisms and evaluate their trade-offs for multihomed enterprises, in order to derive "best common practices" for multihoming. Their most important results shows that there is no benefit to using historical samples, route selection should always be based on current samples. Furthermore they even find that that at times employing historical data to make decisions about which path to forward on can prove detrimental.

2.3 Determinism in Computer Networking

2.3.1 Deterministic Networking Specification

Whilst Time Sensitive Networking (TSN) is able to provide the desired reliability and timing that critical applications require, there is a need for this sort of support across larger networks which can no longer be classified as Local Area Networks (LANs) [13, 18]. To this extent, Deterministic Networking aims to enable the migration of applications suited for TSNs to larger packet switched networks, while maintaining support for existing packet network

applications over the same physical network. [13].

In the context of this thesis, "determinism" does not refer to compliance to the IEFT DetNet Working Group's definition of a Deterministic Network (DetNet). However it is still important to consider their specifications and even more so, their definition of a Deterministic Network, from which the definition of deterministic backhaul can be built.

The Deterministic Networking Problem Statement contains certain requirements of the network as a whole. This is not possible for this thesis, which addresses backhaul over resources not controlled by the operator. The WAN Connector, once implemented, will only be able to select from the available backhaul connections, it has no control over their network. Beyond this however, the problem statement's definition of determinism can be used. This entails absolute guarantees of minimum and maximum latency as well as packet loss and jitter, and secondly the use of more than half of the network's bandwidth. The latter is a crucial aspect of the definition since it bans massive over-provisioning of resources, which is an inefficient solution.

2.3.2 Traffic Shaping, Fair Queuing, and Bufferbloat

Since the DetNet Architecture specification [14] explicitly mentions the use of a traffic shaper to avoid contention loss but also warns against causing bufferbloat in trying to achieve this, it is worth exploring what these are and what they entail.

First, to bufferbloat; this phenomenon refers to the use of disproportionately sized buffers, and their effect on the latencies experienced when these buffers are full [8, 16]. The TCP congestion avoidance algorithm depends on dropped packets in order to detect when the bottleneck link along the path is saturated, and reduce its window size in response. When congested links finally do drop packets if these packets are the last ones in the queue it means the whole rest of the queue must be sent first before the dropped packets are detected and the window size is reduced. This means that there is an unnatural delay before congestion is

detected, and during that period all packets passing through the node experience unnecessary latency due to the congestion induced buffering. This problem has been heavily exacerbated by the cheap availability of memory and the desire for low packet loss. This is a fundamental misunderstanding. Dropping TCP packets is a means of communication and is not at all the same as packet loss which occurs to equipment failure or transmission errors. By dropping packets when a TCP flow exceeds its throughput, congestion is communicated. This problem can occur with UDP as well but it is even harder to detect in those instances.

The first solution to this problem lies in properly sized buffers, another solution is Early Congestion Notification (ECN), wherein TCP is explicitly warned of congestion through flags in the IP headers of messages going back to the sender, instead of through dropped packets. The other solution comes in the form of traffic shaping. Specifically, by using active queue management and fair queuing it is possible to selectively drop packets from a buffer (e.g. only those for a greedy TCP flow) and also dequeue other packets earlier (e.g. a different flow which is not congesting the link). This overcomes the bufferbloat problem by causing the greedy TCP flow to reduce its throughput and letting through those flows that are not using more than their share of the bandwidth.

Beyond mitigating bufferbloat, traffic shaping can still be crucial for highly latency sensitive applications. Contention loss and queuing delays can be avoided by smoothing out any large bursts of packets. Furthermore, which traffic shaper to use is a fundamental design decision that must be made. If a link is to be kept busy at all times, so that maximum throughput is achieved, it means that latency critical packets may not always be forwarded on time. Conversely, by reserving some portion of the bandwidth in case latency sensitive traffic arrives, the full bandwidth cannot be used at all times and thus the throughput is reduced. This is a fundamental trade-off which cannot be avoided.

A traffic shaper is also important for reducing packet loss. In packet switched networks the only source of packet loss which does not arise from a lower layer is when packets are dropped due to buffers overflowing. This occurs naturally in networks with "bursty" traffic

because the natural distribution of bursts of packets mean that occasionally there will be overlapping bursts of packets which cause a link to become congested and drop packets. This is where traffic shaping and queue management is crucial. The bursts may be smoothed out before the bottleneck link, by releasing packets at a slower rate or dropping those that can afford being dropped (if there are any!). It should be noted that a similar effect can also be achieved through resource reservation and/or through a scheduler which makes sure to never over-schedule the network's resources. However this requires exact knowledge of the maximum bandwidth requirements of all the packet flows in the network.

Chapter 3

Approach

To discuss the solution implemented for this thesis requires a review of, firstly, the requirements contained in the problem statement, and, secondly, a description of both the decisions made as well as the process behind making those decisions. To that extent, this chapter will discuss, the components needed in order to address the problem statement. It will try to discuss not only what decisions were made, but also, crucially, what decisions were *not* made, and why. The first section will begin with the fundamental nature of the solution's architecture, thereafter a review of the design requirements which arise from the Deterministic Networking specification will follow, finally the precise nature of the solution's internal workings will be presented.

3.1 Requirements

Concretely, a multipath backhaul solution, as has been described in the introduction, will need to be able to transport the 5G traffic from the edge deployment to the core, while respecting the traffic's QoS requirements. This breaks down into a short list of hard requirements.

1. Forward only the appropriate traffic (as is defined in the 5G flow descriptions)
2. Manage multiple outgoing paths
3. Co-ordinate flows and their requirements with the Control Plane
4. Select the best path to enable the flow's requirements and achieve determinism

3.2 Skeleton Structure

At the outset of this thesis it was decided to use a Control User Plane split. The packet processing and forwarding is performed in the user plane, and the control plane makes the high level decisions about which packets to send where, and in this case, on which outgoing interface to send them. The data plane is kept simple and only performs the time critical packet processing, and blindly follows the instructions of the control plane. This type of architecture is very common in modern software-based networks, for example OpenFlow, the Evolved Packet Core (from LTE), and the 5G core all implement this kind of Control-Data-Plane split.

However, this decision immediately locks one into certain positions. For example, this requires a method of communication between the data plane and the control plane, it also means the state machine that represents the WAN Connector becomes more complex, and complexity always brings with it a greater danger for hidden mistakes and fallacies. The benefit is that the user plane function does not need to perform any of its own decisions, thus simplifying its internal architecture, and these components can be deployed at different locations, and scaled up or out with greater ease.

3.3 Architectural Components Required for Deterministic Multipath Backhaul in a 5G Campus Environment

Determinism, in a networking context, refers to the ability of a network to supply specific QoS requirements with extremely high degrees of certainty. IP flows and applications running in these networks are assured they will experience bounded latency, jitter, and packet loss. To this extent, the Deterministic Networking specification (also called DetNet) is a specification with recommendations for how to achieve determinism. However it does not define determinism. A network which does not fully comply to these specifications may still call itself deterministic (but it may not say it is a DetNet). Since this thesis aims to provide deterministic backhaul, it is not required to follow the DetNet specifications in any form. Though, it does make sense to consider what solutions they have proposed for determinism, since ultimately the goals are the same. In this section, therefore, the DetNet specifications will be used as a source of inspiration to glean ideas, but not as an absolute source of truth.

Deterministic networks have been discussed in the previous chapter. In order to guarantee determinism the IETF DetNet working group has proposed an architecture for determinism over IP networks [14]. Their specification identifies four key mechanisms for guaranteeing the determinism of a flow: 1) elimination of contention loss, 2) jitter reduction, 3) service protection, and 4) explicit routes. Since the 5G campus environment may need to use the infrastructure of other operators this rules out the ability to use explicit routes. A key difference between deterministic networks and 5G campus backhaul networks is that the operator may not control all the links between nodes, and that there are only two nodes in the network that are definitely under the administrator's control- the WAN Connectors at the core and the edge. In such a scenario the problem is less like a network and more like a client - server application with multiple paths between the client and the server.

3.3.1 Elimination of Contention Loss and Bufferbloat

In order to prevent packet loss due to buffer overflows, and, just as importantly, to prevent high latencies due to bufferbloat [16] or excessive queuing, contention between flows must be carefully managed. A greedy flow, which uses too much bandwidth, or the presence of too many flows, whose cumulative bandwidth demands exceed the link's capabilities, will have negative effects on all the flows running over the link.

Elimination of contention loss can be achieved by using a traffic shaper and/or rate limiter, and the ingress to any DetNet domain **must** apply such a function. For the WAN Connector therefore a traffic shaper must also be applied to the traffic on ingress, from the RAN, before it is sent across the backhaul links. Since the implementation of a traffic shaper is beyond the scope of this thesis, an existing solution must be used. The traffic control subsystem in the linux kernel (TC) provides several implementations of different algorithms for both rate limiting and traffic shaping. Some of these algorithms have already been discussed in the previous section. For the purposes of this solution several algorithms were considered: Hierarchical Token Bucket (HTB), Hierarchical Fair Service Curve (HFSC) [34], Time Aware Priority Queueing (taprio), Common Applications Kept Enhanced (CAKE) [20], and Fair Queueing with Controlled Delay (FQ_CoDel).

For HTB and HFSC, due to the implementation of these algorithms in the linux tc subsystem, integrating them within the context of this solution would require filters for each flow. This is not necessarily a hindrance yet, however it would introduce additional overhead upon the creation and/or addition of each new flow. Perhaps more importantly, the implementation of the HTB algorithm provides no means for fair queuing. This means that while it is not possible for greedy and/or malicious flows to use up the bandwidth of other flows because HTB guarantees bandwidth, the greedy flows can cause other flows to experience higher latency because it does not guarantee delay. The HFSC algorithm improves on this by offering both bandwidth and delay guarantees. However this makes it far more difficult to

configure, and would require precise calculation of the service curve for each new flow.

Another feasible option available in linux TC would be the time aware priority queuing (taprio). This queuing discipline provides scheduled gates for specific traffic classes. It could be used to schedule flows according to their requirements and/or priorities, and thus provide a form of prioritization, and, to an extent, traffic shaping. Flows can be assigned different windows and congestion due to a greedy flow is avoided if it is scheduled for a different window. However one issue is that taprio fails to provide fair queuing, and thus would need to be configured and possibly adapted as new flows are added and removed, depending on their priority. Greedy flows of the same priority would have to fight each other to use the bandwidth available in their assigned window.

Lastly there are FQ_CoDel and CAKE. Both provide fair queuing as well as active queue management, which means flows are less likely to experience loss due to congestion, particularly congestion caused by one particularly greedy flow. Since the CAKE algorithm was originally designed for routers in WLAN networks, and these are loosely analogous to a 5G campus network (consider the User Equipment as the Stations, and the Base Station as the Access Point), it makes it an attractive choice. CAKE also allows the use of 8 different priority tins, which fits well with the nature of 5G traffic, where different flows have different priorities. Finally, on a practical level, the implementation of CAKE in the linux kernel uses the kernel's "skb_flow_dissector" which exposes a hook point for eBPF ¹ programs. Specifically within a 5G context, this could allow one to attach an eBPF program to allow CAKE to peek inside of GTP tunnels, and differentiate the flows within them. For all the other algorithms mentioned so far, all of the GTP packets would appear to belong to the same flow, and thus none of these algorithms would work at all if the traffic is arriving in a GTP tunnel. This is ultimately the strongest case that can be made for any of the options listed so far. TC does have one other algorithm which hooks into the "skb_flow_dissector", namely the Choose and Keep Scheduler (CHOKe) [26], which does perform active queue management, however does

¹https://www.kernel.org/doc/html/v5.1/networking/bpf_flow_dissector.html#overview

not provide fair queuing. This means that while it manages buffers in order to prevent them from filling up, it does not guarantee different flows a fair share of the bandwidth, which can still lead to flows experiencing increased delay due to other greedy flows.

3.3.2 Jitter Reduction and Latency Guarantees

In the Deterministic Networking Architecture specification it is recommended to adopt time synchronization as well as sending "Time of Execution" fields in the application packets, in order to achieve jitter reduction. For time synchronization, the PTP protocol may serve well, especially since there is an existing implementation- the linuxptp project - which extends it to work over IP networks. As for the time of execution fields, it may be possible to place these in a GTP header, or to combine them with timestamping.

Latency guarantees are not mentioned within the specification but they are important for the 5G campus setting. Specifically when backhaul options may include satellite links it becomes important to consider latencies. Additionally, in a multipath setting, latency can be a useful criteria for selecting between links, and guaranteeing latency becomes easier to do when it is possible to switch links should one of them start to experience greater latency than before.

3.3.3 Service Protection

In order to protect against equipment failure, it may be recommendable to perform packet replication and/or encoding. To address this the DetNet specification speaks of both duplication of flows, and network coding [5]. This can be a clever solution, however Network Coding requires control over and co-ordination between all the nodes in the network (like other parts of the Deterministic Networking specifications). Though to guard against equipment failure and/or packet loss, duplication does provide an option. So too, does forward error correction (FEC). Many FEC schemes work on a continuous stream of bytes, providing correction for

bit errors, as opposed to correcting the loss of entire packets. There are schemes which address this, fountain codes - such as Raptor [24], which encode the data as being made up of multiple discrete symbols, and these symbols can be mapped directly to packets, thus allowing the Raptor FEC algorithm to protect against packet loss as opposed to just bit errors. Unfortunately the open source ecosystem surrounding fountain codes is not as developed as that around other FEC schemes, and, because of the complexity involved in implementing them efficiently, there are relatively few freely available projects or libraries which can be used to encode packets using such a scheme. For reference, ² maintains a list of C/C++ FEC libraries and none of them support the Raptor family of fountain codes (all other fountain codes incur high overhead). Even if it were more feasible to integrate FEC into the backhaul solution it bears questioning how much benefit it brings. FEC is very powerful in situations with consistently lossy links, however paths over the internet tend to experience loss in bursts, as opposed to at a consistent rate. To this extent it may make more sense to duplicate those flows which require a high degree of reliability.

If one does choose to duplicate packets in order to achieve service protection, that means an elimination function is also required, as well as a packet ordering function. This is because multiple packets may arrive on one link before they arrive on the backup path. If a packet got lost on the faster path, it may still show up on the slower one and thus the packet forwarding needs to pause until the missing packet arrives. The DetNet specifications provide both a basic and an advanced Packet Ordering Function (POF) algorithm. This algorithm can be very easily adjusted to also perform the elimination of duplicate packets. The DetNet specification also clarifies how long one should wait to see if the missing packet arrives (or doesn't) - this value "cannot be smaller than the delay difference of the paths used by the flow" and is called the *POFMaxDelay*.

²<https://aff3ct.github.io/fec-libraries.html>

```

/* initialization */
1 foreach flow  $f$  do  $POFLastSent_f \leftarrow 0$ ;
/* Start POF logic */
2 while true do
3   receive packet  $p$ , with sequence number  $seq\_num$  for flow  $f$ ;
4   if  $seq\_num \leq POFLastSent_f + 1$  then
5     /* eliminate duplicates */
6     if  $seq\_num == POFLastSent_f + 1$  then
7       forward  $p$ ;
8        $POFLastSent_f = seq\_num$ ;
9     end
10  else
11    /* in a separate thread */
12    buffer  $p$  until  $seq\_num = POFLastSent_f + 1$  or  $POFMaxDelay_f$  elapses;
13    forward  $p$ ;
14     $POFLastSent_f = seq\_num$ ;
15  end
16 end

```

Algorithm 1: Basic POF Algorithm Adjusted for De-Duplication

The algorithm shown in 1 is the Basic POF algorithm, with a single adjustment made on line 5 to make sure that an incoming packet's sequence number is not lower than the next one we are expecting, in order to properly eliminate duplicates. The algorithm's **if/else** logic in this representation could be unified in a cleaner way, but it was done this way specifically because it thus requires just one extra line (line 5) to be changed from the DetNet specifications POF algorithm, and that one change suffices to make it a Packet Ordering and Elimination Function.

3.3.4 Multipath Considerations

While the previous subsections have all considered requirements for determinism in an IP network, the issue of multipathing still remains. For a component which is backhauling over multiple links in an IP networks this means link and/or path selection is required. Traditionally, routers select links primarily based on their ability to route the packet to the destination, and secondarily based on various metrics, which can be defined by the administrator. For the WAN Connector's scenario routing considerations are not made- it's purpose is to backhaul the traffic from the RAN / Edge to the Core, where the network's own internet gateway can perform the routing.

At this point it is crucial to differentiate, again, between multihoming and multipathing. Multihomed hosts are able to select one of many links for their outgoing traffic's destination. In multipathing the same link may lead to multiple paths to *different* destinations. When two multihomed hosts communicate with eachother multiple different paths are opened up. Multipathing over the internet requires complex data collection about all the routers in between the source and the destination, see [15] for an example of such an implementation. This is because multipathing over the internet needs to choose paths with maximally edge-disjoint paths to achieve the best reliability, as well as maximally vertex-disjoint paths in order to reduce the load across all nodes on the path.

The goal of this thesis is of course to use the link and/or path selection to provide determinism for the flows being backhauled. To this extent then, the link selection algorithm must attempt to forward flows over paths which can provide at worst the maximum allowed latency and jitter, as well as meeting the minimum reliability. Here it is worth noting that while duplication cannot really be used to guarantee latency, it can be used to improve reliability since a flow which is duplicated across two paths is far less likely to experience packet loss.

3.3.5 Summary of Necessary Features

Looking back on the previous subsections, the following points (in no particular order of priority) were identified as mandatory for a deterministic backhaul solution: 1) traffic shaping, 2) path selection according to jitter and latency requirements, 3) service protection (i.e. via packet duplication), 4) packet ordering and de-duplication, and 5) time synchronization. The ways in which these can be addressed, or quite simply the way in which they were implemented, will be discussed in the following section

3.4 Overview of the WAN Connector's Features and Components

At this point, one can take the skeleton structure proposed in the first section, and superimpose the other requirements which were determined in the second section. Combining these ideas yielded the final implementation, which will be now be presented.

3.4.1 Path Selection Algorithm

Meeting the various requirements - jitter, latency, and delay - of a flow can be formulated as a multi-constrained QoS problem. Solving such a multi-constrained QoS problem via path selection is a binary optimization problem, and the key to deterministically meeting the flows requirements. The problem can be posed as: "select those paths on which to forward packets while making sure to satisfy the latency, jitter and reliability requirements of the given flow, and minimizing the overall weight of the paths used". The mathematical definition is as follows:

$$\text{Minimize } \sum_{i=1}^P w(x_i) \quad (3.1)$$

$$\text{Where, } d(i) * x_i \leq D \quad (3.2)$$

$$j(i) * x_i \leq J \quad (3.3)$$

$$1 - \prod_{i=1}^P (1 - r(i) * x_i) \geq R \quad (3.4)$$

$$\text{for } x_i \in \{0, 1\} \quad (3.5)$$

Here the variables D , J , and R are the flow's delay, jitter, and reliability requirements, while the functions $d(i)$, $j(i)$, $r(i)$, $w(i)$ are the estimated delay, jitter, reliability, and weight of link i . The predicted values will usually just be the latest measurement, as recommended in [7], however there is room here to use more advanced metrics to predict the future link quality and thus perform preemptive path switching in future work. The total number of paths is P . The x_i variable indicates whether or not link i shall be used. If a solution is found, then the flow's packets will be forwarded on each link i where $x_i = 1$, and if no solution can be found which satisfies these conditions then the flow is rejected because its QoS cannot be guaranteed.

It is worth noting that solving such problems is NP-Hard [23]. However this hardness arises primarily because of equation 3.4, the equation for reliability (also the only non-linear equation). Due to this equation one must consider every possible combination of paths on which to forward, and the complexity is $O(2^n)$. This means, even though Akella et. al [6] have shown that multihomed approaches experience diminishing returns after more than 4 links, attempting to brute force the solution by limiting the number of outgoing interfaces to 4 still yields a very large problem space - in the worst case both WAN connectors could have 4 outgoing paths, leading to 16 possible paths and thus 65536 possible combinations to consider.

In order to further avoid the combinatorial explosion the problem needs to be parameterized even more. The first logical parameterization has already taken place by limiting the number of interfaces to 4. This can be expanded on by limiting duplication to only take place over disjoint interfaces. The reasoning behind this is that in a geographically distributed Campus 5G environment, and especially for environments featuring wireless backhaul (e.g. satellite), it is more likely that a link's reliability is most affected by the over the air transmission on the first hop. By performing this parameterization the problem space shrinks considerably, to just 2^4 possible combinations of outgoing paths, but, because certain paths going out on the same interface are ignored, the optimality of the overall solution is gone. This is an acceptable trade off for quick computation, and reasonably strong guarantees on reliability. When choosing which path to include from a set of paths using the same outgoing link, only the path with the greatest reliability is taken into consideration so that the likelihood of rejecting a potentially viable flow is as low as possible.

3.4.2 Packet Ordering and Elimination Function

Since it is possible for flows to be duplicated across multiple paths, it becomes a necessity to have a packet ordering and elimination function. To be able to re-order and de-duplicate packets means that their sequence numbers need to be tracked. This thesis' implementation will track the sequence number on a per flow basis, using the GTP "sequence number" field. This thesis uses it's variation of the DetNet POF specification's Basic POF algorithm [38], as presented in 1.

3.4.3 Internal Architecture

Control - Data Plane Communication

A protocol is desired which can provide reliable communication over multiple paths in order to communicate between the data plane and the control plane which is where the statistics are collected, flows are added or removed, and the multipath decisions are made. This is crucial, since in a geographically distributed campus 5G deployment it is unlikely that there will be a separate management or control network which uses different underlying network infrastructure than the data plane. Therefore link failure, as well as packet loss, on the link being used by the control plane must either be avoided or protected against. To this extent there are only two feasible options - either multipath TCP or SCTP. SCTP does support multihoming and intelligent failover, despite not being explicitly designed for multipathing, and is slightly easier to manage and configure than multipath TCP, which is why it was chosen for the data plane control plane interactions.

GTP Tunneling and Custom Header Extension

The tunneling protocol used between the two data plane instances of the WAN Connector will be GTP version 1 [2]. The specification for GTP allows for the use of sequence numbers, as well as the use of extension headers. For the data plane communication there will always be an additional custom header sent which includes a timestamp taken by the sender. The timestamp is taken in microseconds. For the timestamp there are several options available within linux. The TAI clock in linux (representing the Atomic International Time) is used in this implementation because it does not have leap seconds, and so is a monotonic function, which is an important guarantee for time sensitive applications. An example of a packet capture of the WAN Connector using the custom header is included in the appendix C.

Normally, in 64 bit linux systems, a timestamp takes up 16 bytes and consists of 8 bytes for

the time in seconds since the epoch, and another 8 bytes for the nanoseconds. To reduce the footprint of the timestamp in the header, the seconds are represented with just 8 bytes, thus wrapping around the interval $[0, 255]$. The receiver needs to take this into account, but is programmed to do so, and since the path delays are not realistically expected to ever exceed 3 seconds this is more than enough. The microseconds have a maximum value of 10^6 and can be represented with just 24 bits. For Wide Area Networks and internet connections usually the latencies are on the order of milliseconds and as such microsecond precision is deemed to be sufficient.

Metric Collection

To be able to intelligently switch flows between paths, and to be able to know when this is required, it is necessary to collect the relevant metrics about latency, jitter, bandwidth usage, as well as packet loss. Packet loss can be detected via the sequence numbers, while delay and jitter can be calculated using the timestamps passed along in the GTP headers. These metrics are periodically reported to the control plane so that it can make its decisions on up to date data. Keeping a healthy overview over the state of each path requires periodic probing on these paths. This is especially important for detecting when paths become viable again, since these paths will not have any traffic on them while they are considered down, or if they have experienced high latency and/or jitter recently. The probes are sent once per period of reporting so that they have a minimal impact on the bandwidth usage.

Flow Descriptions, Flows, and Hashing

Incoming packets need to be quickly mapped to their respective flows. It is common in network environments to perform flow hashing, in order to quickly lookup which flow incoming packets belong to. This approach makes sense here too. However, since 5G flow descriptions can apply to various IP and port ranges as well as matching based on the transport layer

protocol it is not possible to hash an incoming flow and obtain the same hash as the flow description. Furthermore it is possible for a flow to match multiple different flow descriptions. In these cases the first matching flow description is taken. This implies some sort of ordered storage of flow descriptions will be required, as well as a method to quickly match packets to flows. The solution used here is to maintain a simple linked list of the known flow descriptions, and match packets to their flow descriptions, if the packet is unknown. For "known" packets, which have been matched to a flow description, these are hashed and stored in a table, alongside their flow descriptor and a list of paths on which the flow is supposed to be forwarded. This allows quick lookup for every subsequent packet, as well as making it simple to change which paths flows are meant to be forwarded on. One alternative to storing the flow descriptions in a linked list would be to store them in a tree, such as in [37], sorted either by priority or ID or even by the hash value of the flow description. This method would provide a faster lookup of new, "unknown" packets, whose hashes don't yet match to a flow, but since new flows are not such a frequent event the additional complexity of implementing such a data structure may not be worth the gain in performance.

It is important to make a good choice when selecting a hash function. Since the hashes used are only for the purpose of lookups it would not be recommendable to pick a hashing algorithm designed for use in cryptography. These algorithms are designed so that it is very difficult to reverse engineer the original value from the hash, and this may often make the computation of the hash more computationally intensive than when computing hash algorithms designed for looking up table entries. Lastly, for hashing algorithms it is desirable that they provide a healthy distribution, to reduce collisions. Collision reduction can also be affected by choice of hash table size- and in general it is usually recommended to use a prime number. In this implementation, the MurmurHash algorithm was chosen due to its strong performance for lookup-based hashing. Specifically, the MurmurHash3 [10] version was chosen, and since it can generate 32 or 128 bit values - the exact sizes of IPv4 and IPv6 addresses - this simplifies the hashing implementation for IP flows. When a packet comes into the WAN connector it is hashed based on it's destination address, source address, destination port, source port,

and finally the transport layer protocol number. If the packet was tunneled with a GTP header, the header is removed and the hash is performed on the tunneled packet, not on the GTP packet. If the hash does not match, the packets is compared to the list of existing flow descriptions. If the flow belongs to one of these descriptions then that decision is stored in the hash table.

3.4.4 Overview

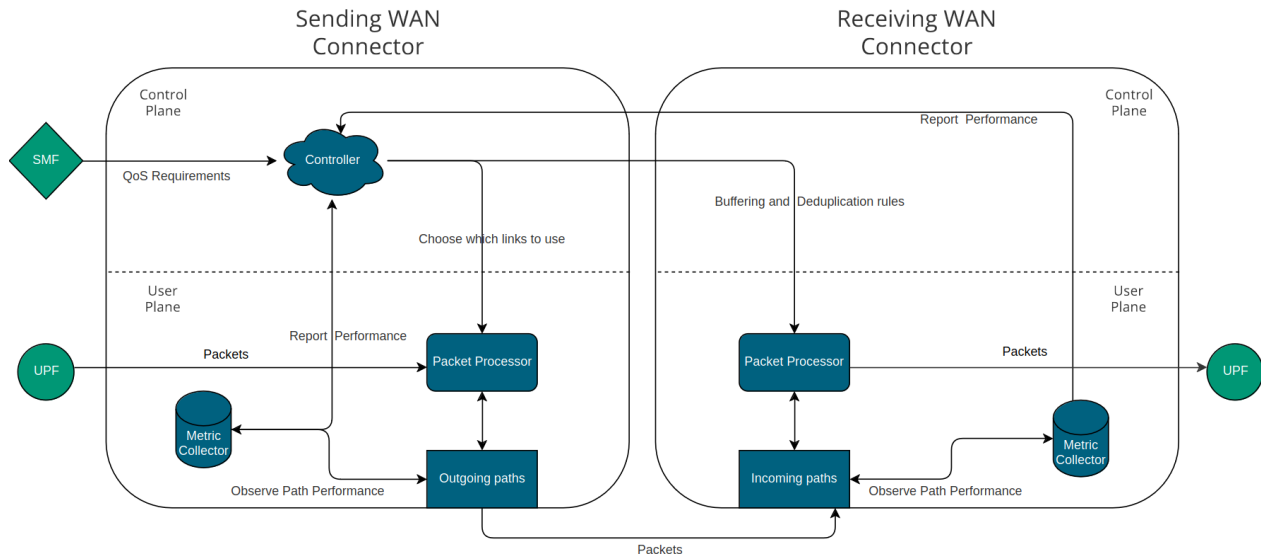


Figure 3.1: Internal Architecture of the WAN Connector

The internal architecture of the WAN connector is shown in Figure 3.1. The connector is split into a control plane, which uses the algorithm in 1 to choose which interfaces to forward on, and a data plane, which performs the forwarding, as well as any necessary buffering, re-ordering and/or de-duplication. The WAN connector also performs different tasks depending on whether it is the origin or termination of a flow. For example, the terminating node may be receiving duplicate packets on the other paths, and it must know to drop these. Additionally, the sender, can only store statistics about how many bytes it has sent, and on what paths. The receiver is able to use the information contained in the packets it has received to construct a complete picture about the nature of the path - its jitter, latency and

also packet loss. The link selector in the control plane can then use this information to make its decisions. When deploying the WAN Connector there is only one Control Plane unit (the controller), and two Data Plane instances. This is because it would be redundant to deploy an additional controller for the remote data plane instance. In the event that each link fails then the remote controller cannot help, even though it can detect it, and if at least one link remains up then the core controller will switch to using that link and can inform the data plane unit to do the same with the traffic it is forwarding.

One part which is missing from this diagram is the traffic shaper. That is because it is not connected to the control and the data plane, rather it is meant to be installed alongside them. Furthermore, it is important that the traffic shaper resides in the part of the network before the packets enter the WAN connector. If the deployment features a UPF in the edge network, it may make sense to place the traffic shaper before this as well. For the core, the traffic shaper will similarly need to be placed at the ingress point at which packets reach the WAN connector.

Chapter 4

Evaluation

This chapter reviews how the approach from the previous chapter was tested, and discusses the observed results. This chapter will explain the setup of the testbed, as well as the motivation behind the types of tests which were performed. Finally a discussion of the results and their origins will take place.

4.1 Approach to Evaluating Performance

In order to evaluate the success of the proposed approach a multipath scenario with four links will be emulated and investigated. Each test will consist of two WAN connectors, and the four emulated links going between them. The characteristics of the links will be changed during the tests, each time with respect to different properties: latency, packet loss, and jitter. The evaluations will be performed on a testbed consisted of four separate host servers. Two of these servers will act as the originating and terminating WAN Connectors. One server will generate and measure traffic, in between the two WAN Connectors, a server will act as the link emulator. This architecture is shown in figure [4.1](#)

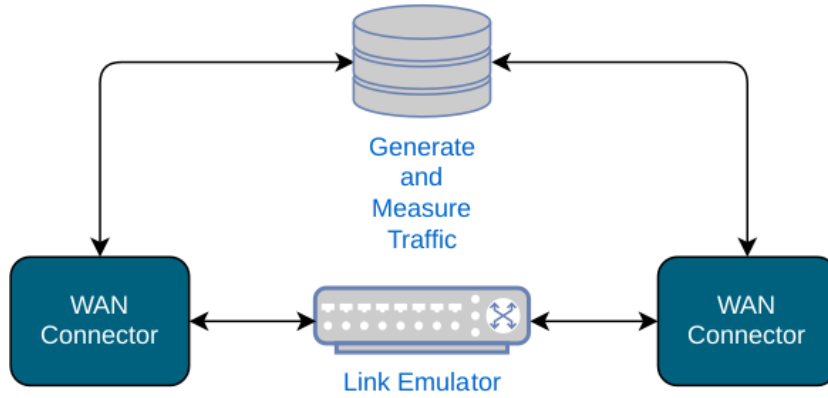


Figure 4.1: Testbed Setup

In figure 4.1, the links between the servers are 10 GB ethernet cables. The hosts are connected as is shown. In the "traffic generator and measurer" two linux namespaces are used to separate the traffic generation and the measurement. The packets generated in one namespace are routed via the ethernet connection to the next host (the terminating or "edge" WAN Connector), because of the linux namespace they are not aware that their destination is technically the same host. The packets are then backhauled by the WAN Connector over the emulated links in the next host, and via the terminating WAN Connector, back to the original host, but now in the "measurer" namespace.

The link emulation is done using the tc subsystem. First a root Hierarchical Token Bucket qdisc is established on the "downlink" and "uplink" interfaces (the ethernet cables between the WAN Connectors). In this case, the link emulation is performed before the outgoing packet is queued on the ethernet cable, so that is also where the qdisc is placed. Secondly, the qdisc is given children classes, one class per link. This allows one to establish bandwidth caps for each emulated interface. Within these HTB classes the root qdisc is a netem qdisc. These allow the installation of link characteristics such as latency, packet loss, and jitter.

In order to make the emulation more realistic, the netem qdiscs will periodically be adjusted, this allows one to degrade or improve links over time, for example by increasing the latency or packet loss in frequent intervals, up to a large value. It also more closely mimics the real-life behavior of WAN connections, which do experience changes in their characteristics

over time.

For the purposes of evaluating the WAN connector, three series of experiments will be performed to isolate and investigate its link switching capabilities. First, purely latency based link selection will be investigated- a flow will be defined with specific latency requirements and then the emulated links will have their latency repeatedly changed. The same will be done once with packet loss and with jitter. In each of these scenarios the other two parameters will be held constant, so that the WAN Connector can be judged based on its ability to meet one of the latency, jitter, or packet loss requirements, without the other two requirements interfering. In order to simulate realistic scenarios these experiments will also be repeated once with background traffic which aims to saturate the link's capabilities.

The data for the emulation used in all of these scenarios was obtained using the Satellite Constellation Network Emulator (SCNE) ¹ from the European Space Agency. The default scenario is 24 hours long and consists of 288 satellites, 10 gateways (GW), and 50 user terminals (UT). The reason this data is used is because LEO satellite links are expected to be the most common type of backhaul link for geographically distributed 5G Campus Network deployments. Furthermore, since LEO satellite connections are exceptionally more volatile than "regular" links [12] [25]. By choosing to emulate more volatile links in the testbed, the WAN Connector experiences a more difficult test, and one can be assured that the real life performance in "normal" scenarios, is highly unlikely to be worse.

Since the WAN Connector is defined to work on at most four outgoing links, each investigated scenario will feature four links. These links will all be based on data from the SCNE emulation.

¹<https://connectivity.esa.int/projects/scne>

4.1.1 Accuracy of the Emulation

To make sure that the testbed emulation also matches the values from the SCNE, a brief comparison was done, where the testbed's latencies were compared with those recorded in the SCNE.

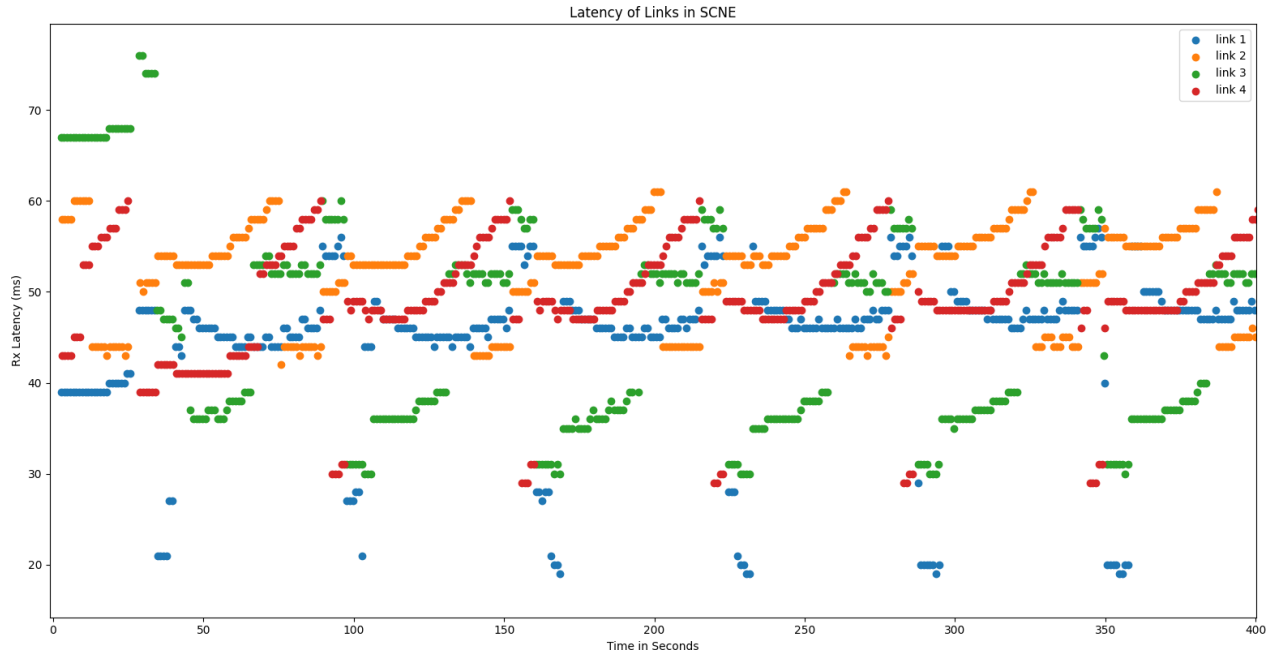
Figure 4.2 shows the comparison between the data observed in the testbed, and the SCNE. The testbed data in the graphic is measured by the packet flows running over the individual links. The emulation data comes from the values recorded during withh the SCNE. Since the testbed data is collected at a rate of 100 packets per second, but runs at 10 times the speed of the SCNE emulation, there are 10 times as many data points. This explains why the testbed data appears blurrier.

As can be seen in the figure ??, the latencies recorded on the testbed closely match those from the SCNE, which it was seeking to replicate.

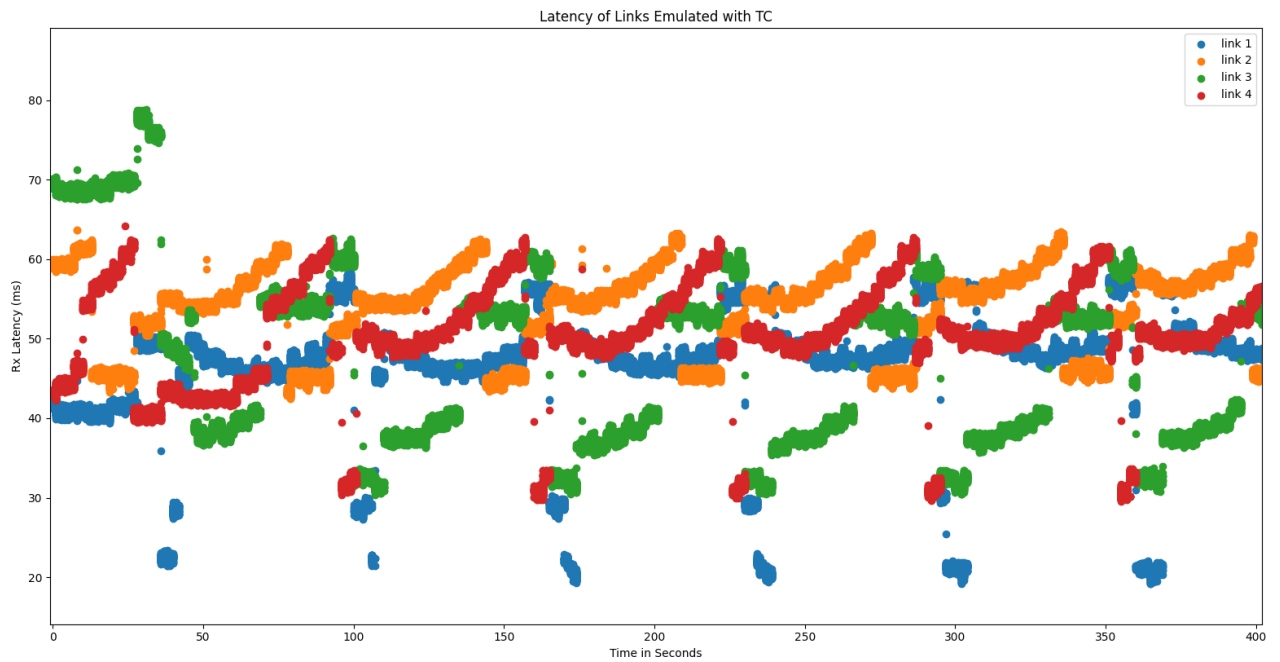
4.2 Latency Based Path Switching

First is the scenario investigating latency. During the emulation, the simulation data is used to adjust the tc netem qdisc on the "link emulator" machine periodically. Because the SCNE simulator provides latency data over a 24 hour period, in 10 second intervals, the simulation data is "sped up". That is the emulation is changed each second, based on the next datapoint from the simulation. This means the 24 hour window can be run over a 2 hours test. Since real links are usually steadier and less volatile than the sped up emulation one can expect results in a real situation to be better.

In the experiment 100 packets are sent per second across each link, as well as 100 packets per second from the traffic generator, through the WAN Connector. This is done to have a comparison between the performance of the WAN Connector and the baseline performance



(a) Latencies Recorded in SCNE



(b) Emulated Latencies Measured in the Testbed

Figure 4.2: Comparison of the Emulators Latencies and the Testbed

level, which would be the best performing link. In any scenario the minimum acceptable performance of a path switching application would be to provide a better performance than the best single link. Otherwise it would be better to just select the best link, and not use the application at all.

4.2.1 Performance Compared to Single Links

For the experiment a flow is defined in the WAN Connector with a minimum round trip time (RTT) of 56 milliseconds. This value is chosen because it is the mean of the mean values of the latencies of the four links from the emulation. For the evaluation of the latency based path switching a single flow matching the definition stored in the WAN Connector is backhauled through the WAN Connector at a rate of 100 packets per second. In order to serve as a comparison, identical 100 packet per second flows are sent across the four other links. Finally their latencies are compared, and the latencies measured by the four links are used to calculate the Oracle approach.

The graphics [4.3](#) and [4.4](#) show a cumulative distribution of the different latencies experienced by the packet flows running over the individual links and the WAN Connector, as well as the performance of the Oracle approach. The Oracle results are purely theoretical. For each packet, it is able to select the best link to forward on, based on that links future characteristics. This approach is calculated by comparing each packet sent by the WAN Connector with the latency that packet *would* experience on any of the other links, and if the latency of the WAN Connector's link would be greater than the minimum value (56 ms), selecting a different link (if one exists) with a lower latency. Since it always picks an optimal path, the Oracle approach provides an upper bound on the best possible performance. In practice it is not possible because it has knowledge of the future, but nonetheless it acts as a very valuable comparison.

The CDF shown in figure [4.3](#) and zoomed in on in figure [4.4](#) tracks the probability that the

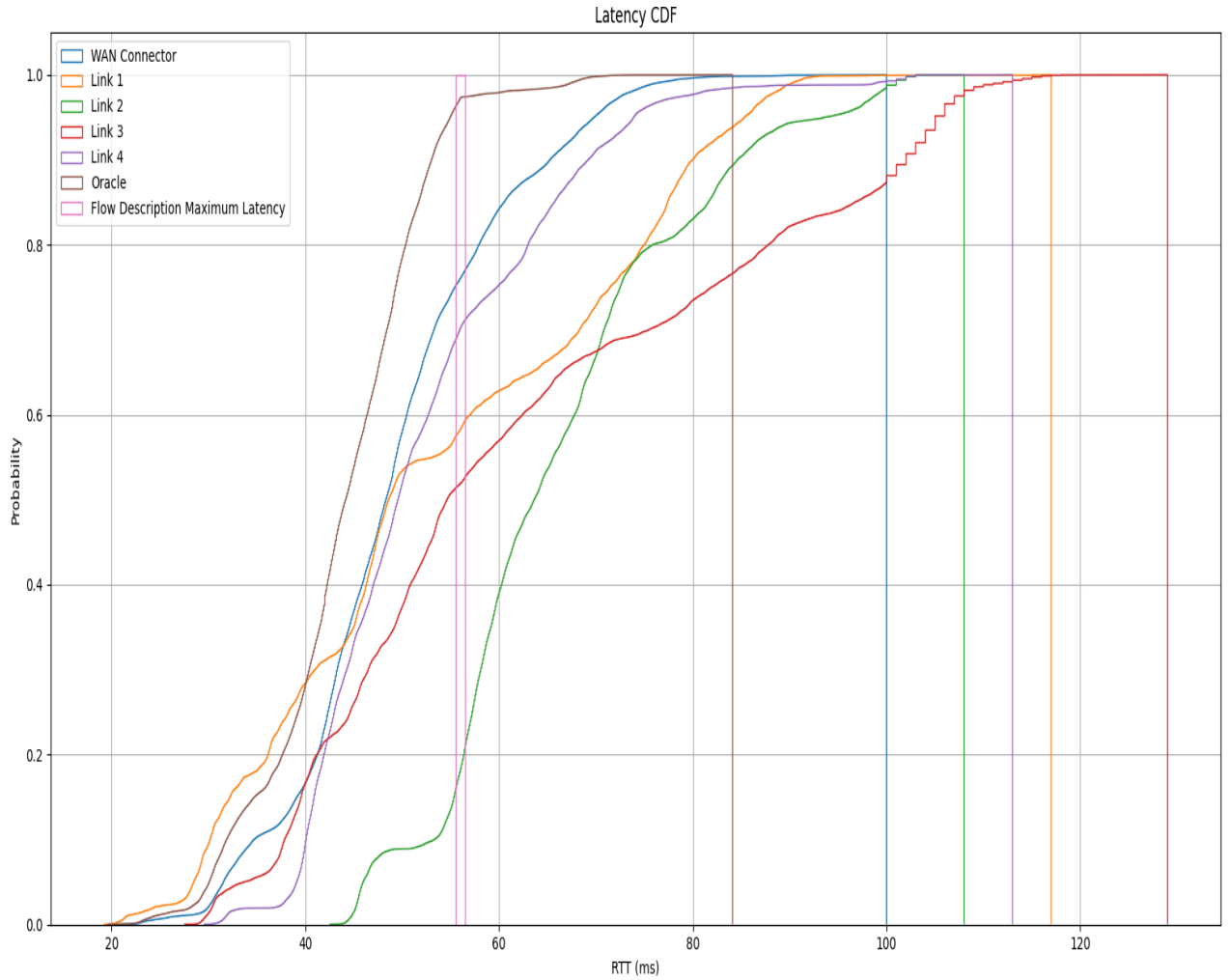


Figure 4.3: Cumulative Distribution of the Latencies

flow experiences a latency less than or equal to the latency value on the y axis. The pink bar labelled "Flow Description Maximum Latency" is set at 56 milliseconds and represents latency bound for the flow being measured, which was communicated to the WAN Connector.

The results in the figures show that the WAN Connector is able to achieve a superior performance to any of the individual links on their own. This was the baseline requirement and it is important that it is able to achieve this. However the performance does not significantly improve on the next best link, which is is able to achieve the maximum allowed latency 72% of the time, while the WAN Connector does so 77.5% of the time. The Oracle

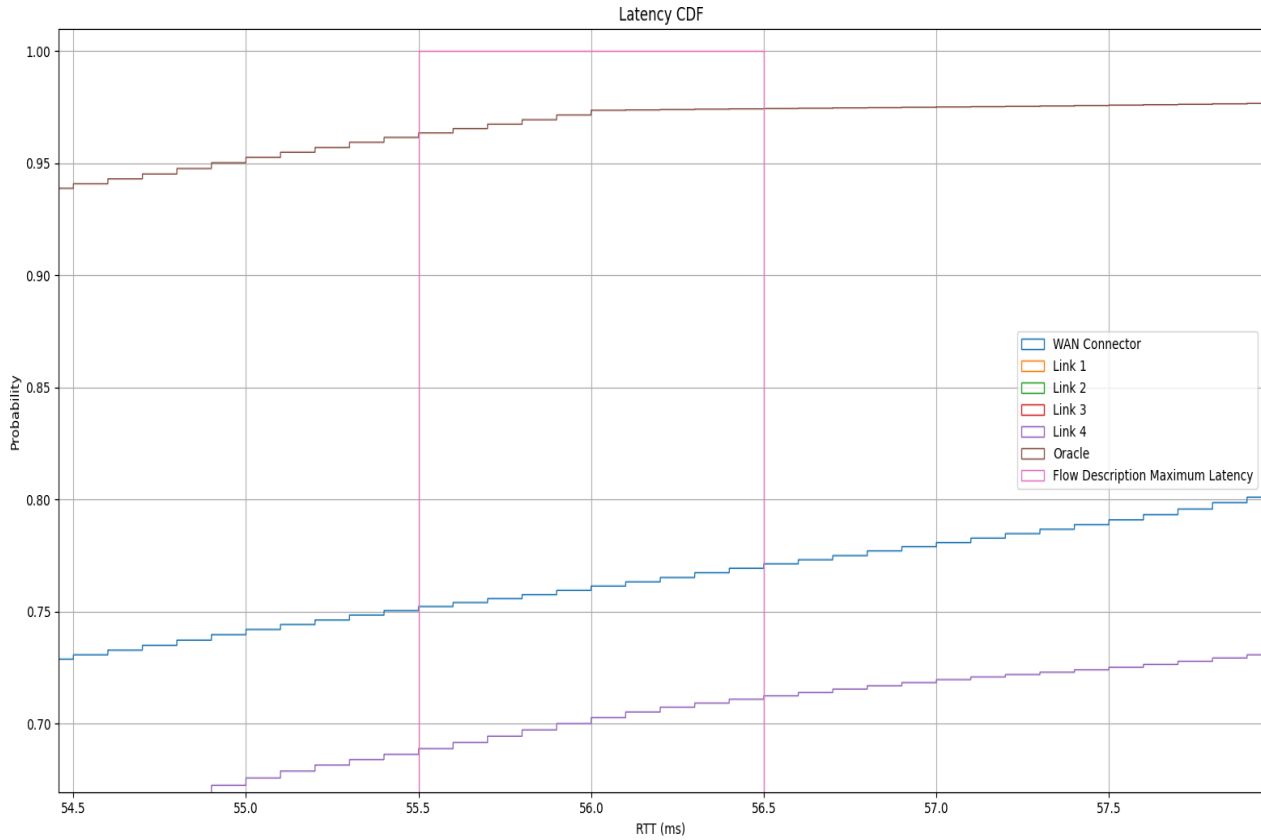


Figure 4.4: Closeup of the Latency CDF

approach shows that in theory one could, at best, have maintained the minimum latency required by the flow in 95% of the cases. This means the WAN Connector’s performance is only 80% optimal.

4.2.2 Latency Performance Under Load and the Importance of Traffic Shaping

During this experiment, the same 100 packet per second flow from before is repeated, but background traffic is added. An iperf3 test is run across the WAN Connector at the same time as the latency critical flow, and the background traffic is given the same latency requirements, so that the WAN Connector always schedules both flows to the same link. This ensures that

whatever link is being used will be saturated. Latency performance during congestion and close to congested scenarios is a crucial element of a deterministic backhaul solution. The WAN Connector's control plane only makes decision about which link to forward on- it does not perform load balancing or prioritization. The data plane does not do this either, it implements a purely First In First Out (FIFO) approach to packet queuing. This is done because the burden of implementing an appropriate packet shaper, in addition to the other multipath calculations, is too significant for the purposes of this thesis, and because very powerful traffic shaping implementations, which are easy to integrate into this solution, already exist.

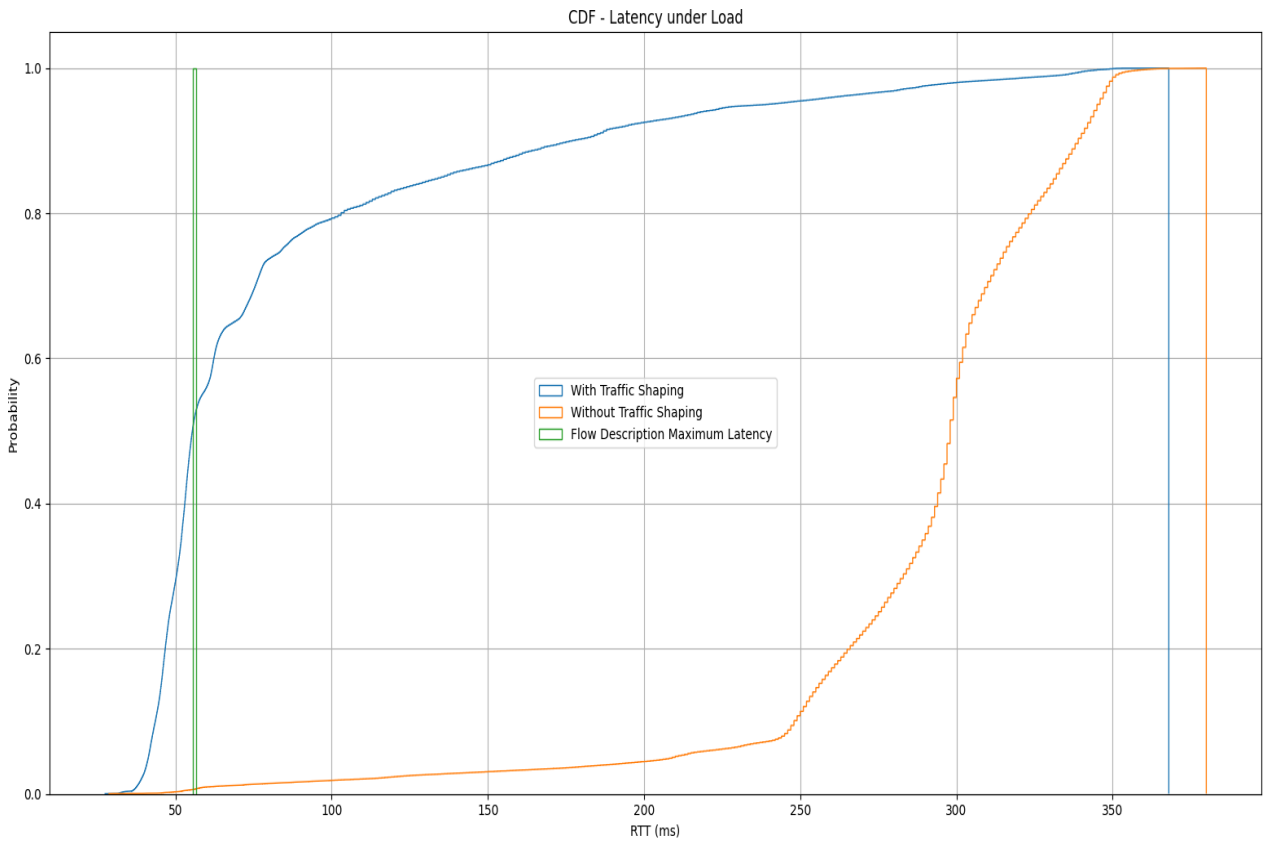


Figure 4.5: Latency with Background Traffic

The figure 4.5 neatly demonstrates why traffic shaping is an essential element of any deterministic network. In the case in which there is no traffic shaper the latency critical flow experiences significant additional latency due to the network being overloaded. The TCP algorithm fills

the link's capacity, and the latency critical packets spend too much time buffering and arrive late. It should be noted that the WAN Connector's performance with the traffic shaper is worse than in the previous scenario, where there was no background traffic, and only 50% of the packets are under the required latency of the critical flow. However the approach without the traffic shaping is at just 1%.

4.3 Reliability Based Path Switching

In the next scenario, to analyze the ability of the WAN Connector to achieve the required reliability of a given flow, the same test flows from before are run over the WAN Connector and the four outgoing links. However this time instead of varying the latency, the packet loss ratio of the link is changed. The SCNE emulation's data is used for the testbed once again. However since the emulation does not provide packet loss on a per second basis, the latencies of all the different links in the simulation are used. For the emulation, these 51 values are looped over in 10 second increments. Each link is given a shuffled set of these 51 values. This means that the cumulative value of the packet loss experienced by all of the links is the same (2%), but the instantaneous values will differ. To appropriately analyze this scenario the flow being backhauled over the WAN Connector is given a reliability threshold of 0.1% this is because in this scenario, where the existing links all provide an average of 2% reliability, it becomes necessary to duplicate at some points, in order to achieve the required reliability. Afterward, the experiment is repeated under load, like in the latency scenario, where the test flows across the links themselves are not run, and a second iperf3 flow is run across the WAN Connector in order to saturate the link. Like before, this second flow is given the exact same requirements as the test flow, in order to ensure that they are always running across the same link.

The results of this experiment are shown in figure 4.6. The packet loss recorded by each of the flows over the period of the experiment is plotted in a bar graph, as well as the theoretical

packet loss which would have been experienced by the Oracle approach.

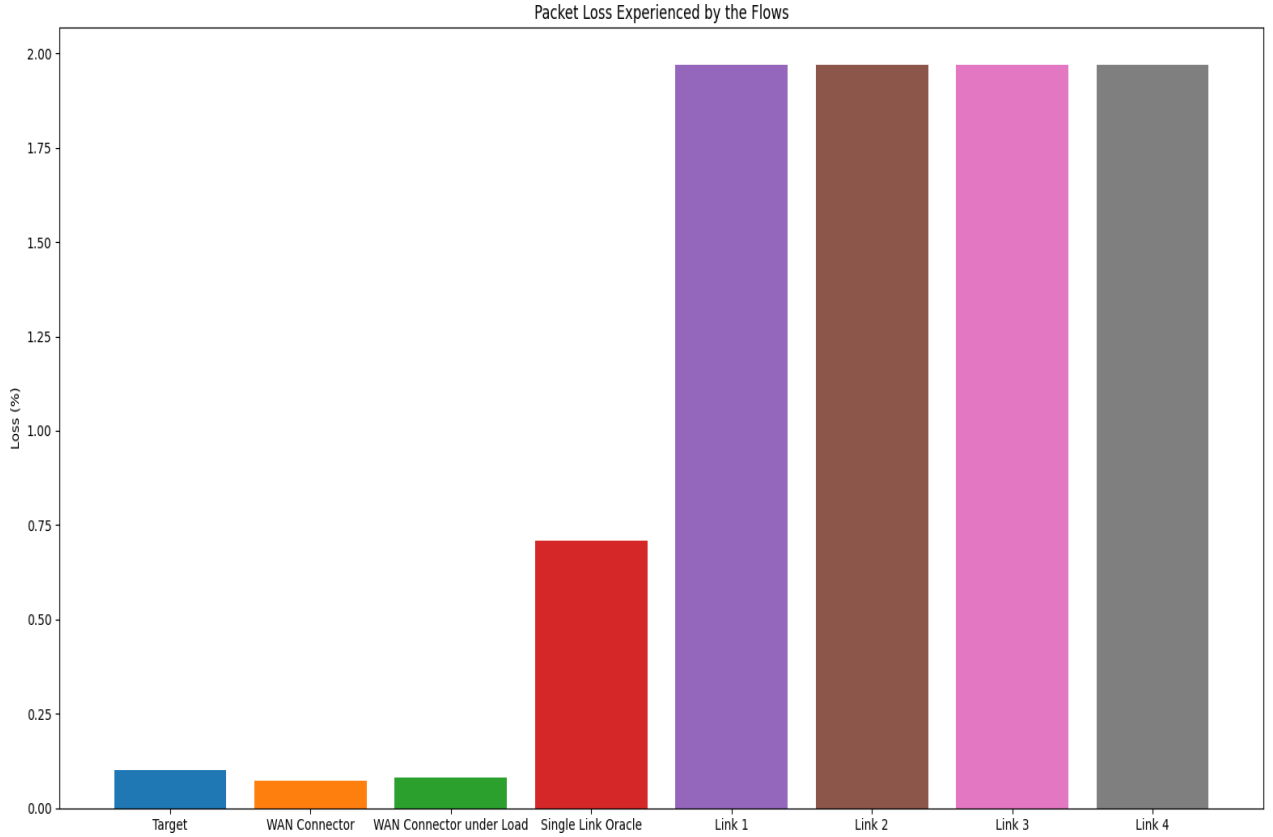


Figure 4.6: Reliability

For the analysis of the performance this time the Oracle approach is defined as choosing, before forwarding any packet, the link with the lowest packet loss ratio. But, crucially, the Oracle is not allowed to duplicate flows across different links. This explains why in figure 4.6 the Oracle is outperformed by the WAN Connector. Indeed, as the graphic shows, it is not possible to achieve the critical flow's required reliability with any of the available links, nor with the Oracle approach. But the WAN Connector is able to achieve it. This speaks to the potency of the flow duplication approach, which is made possible in the optimization equation used to decided on which flows to forward.

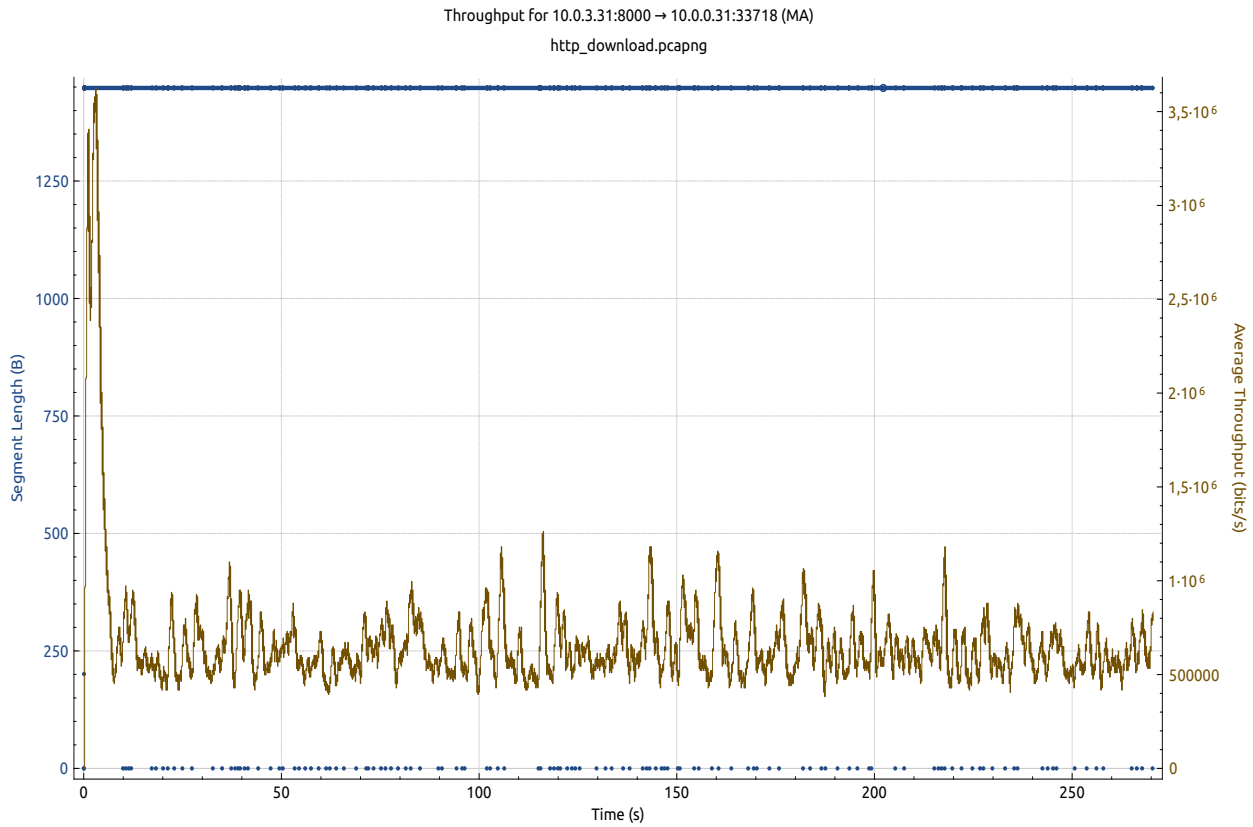


Figure 4.7: Throughput of a Duplicated Flow

4.3.1 Impact on Throughput

One of the unfortunate impacts of packet duplication is that it reduces the available bandwidth. It represents a waste of up to 50% of the available resources, in the worst case, where every flow must be duplicated. Additionally, since flows which are being duplicated need to be re-ordered, this means that the FIFO approach for packets cannot be used and instead the packets must pass through a Packet Ordering and Elimination Function (POEF). When a packet arrives out of order the POEF starts to incur significant additional overhead because it must store all packets until the missing packet arrives, or a point is reached at which the stored packets must be dequeued. This point is either the point at which the packets expire, or the queue of out of order packets becomes full. Beyond this, the storing and eventual mass de-queuing of out of order packets can also interfere with the TCP algorithm- thus leading

to reduced throughput.

These two effects - the POEF overhead and its effects on TCP- lead to a very heavily reduced throughput, as can be seen in figure ?? . The graph in the figure shows the throughput experienced by a routine TCP download, running over the WAN Connector, and experiencing replication, as well as the segment length of the TCP stream. As can be seen in the graph, the throughput is very low, and is not able to saturate the link's bandwidth despite being the only flow running at the time.

4.4 Jitter Based Path Switching

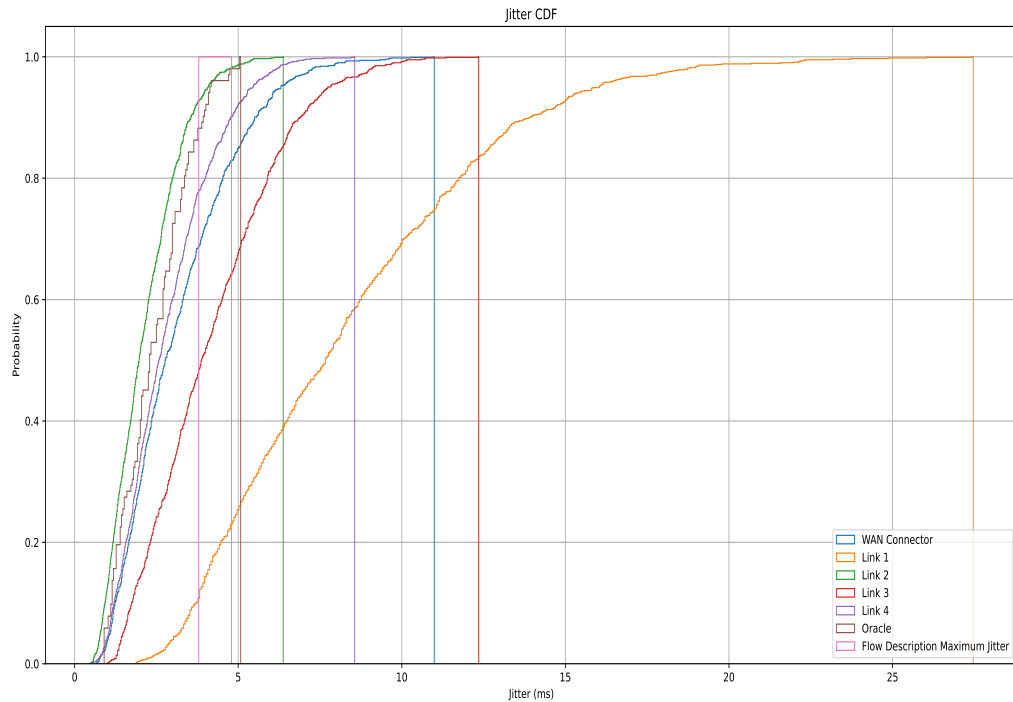


Figure 4.8: Performance of Jitter Based Path Switching

For the final analysis, the jitter must be investigated. For this purpose the simulation data

was once again adapted. This time the standard deviation of the latency experienced was calculated for each of the 102 connections in the simulation. Then these standard deviations were scaled down by dividing them by one, two, three, and four, respectively, in order to acquire four different set of jitter characteristics. Finally the scaled values were shuffled. This resulted in 4 links with varying jitter characteristics, which suffices for the emulation. Each different jitter value is emulated for 10 seconds on the respective link. During the emulation scenario a UDP iperf3 stream is ran at a rate of 3/4ths of the link's capability for 1000 seconds. iperf3 records jitter each second, and this data is used to then analyse the performance. Since more than half of the link capacity is being used by the flows being sent on the links (and not through the WAN Connector), the WAN Connector cannot be tested at the same time. Thus, to test the WAN Connector the emulation is repeated but with the flow only being sent over the WAN Connector. Lastly, the Oracle approach is calculated based on the simulation data. The Oracle approach, like before, displays the theoretical optimal performance. It chooses, at each point in time, the link with the lowest jitter. The cumulative distribution of the jitter experienced by the flows is plotted in figure 4.8.

As can be seen from the figure, the WAN Connector gave a very disappointing performance. It displays an inferior performance to two of the links. This means the path switching is actually causing a worse performance. This is most likely down to insufficient statistics causing poor path selection, because jitter statistics can't be gathered about other paths in this single flow scenario. The probe packets which are periodically sent out do not provide enough data points to accurately determine jitter. While they can give a sufficient baseline for latency calculations, as well as helping to determine if a link is dead or alive, their jitter measurements are not likely to provide a strong insight into the true nature of the link's jitter.

A final note about the plot in figure 4.8. In the figure the first link actually performs better than the Oracle for most jitter values, before ultimately proving to have the greater maximum jitter values. This occurs because the Oracle is based on the simulation data and not the

measurements. Therefore it is likely that the better performing link experienced, for the most part, very low jitter, and thus because there are more data points available than for the Oracle (and because no doubt there was some small amount of luck with the random values used by tc's netem qdisc), the percentage of packets below a certain level of jitter is greater. However ultimately the maximum values experienced by the flow, during those times when the link experiences it's worst jitter, are worse than the Oracle approach's worst cases, as one would expect.

4.5 Discussion

In light of these results one can conclude that the approach presented in the previous chapter is not without its flaws, however the results have also highlighted some of the functional aspects that were important for good performance. For latency based path switching it was able to outperform the other available links and achieve 80% optimality, and the inclusion of a traffic shaper proved vital for preventing the latency critical flow from experiencing higher buffering times in the presence of a competing flow. For reliability, the WAN Connector showed that it was able to guarantee a critical flow a lower level of packet loss than would have been possible using any single link. Furthermore, it was even able to outperform the single link Oracle, and in a situation where the single link Oracle is unable to provide the required minimum level of packet loss for the desired flow, the WAN Connector was able to deliver the desired reliability. This is a hugely important result because it shows that the WAN Connector can provide greater reliability than even an optimal single link path selection algorithm, via the packet replication. This can enable critical applications to be run in environments where the backhaul options would usually be too unreliable to enable the critical applications.

This section, however, also demonstrates the downside of the packet replication approach, as it greatly reduces the overall throughput. Beyond this, during the jitter experiments, it is

seen that the WAN Connector provides a poorer performance than simply using a single link for the entire experiment. This is fundamentally a bad result. The baseline by which any path selection algorithm should be judged is its ability to outperform the individual paths, since it can always select a superior option when one exists. Since the WAN Connector fails in this instance it must be concluded that it cannot provide an overall reduction in jitter for any jitter sensitive applications. The author would suggest, however, that the problem in this respect is not due to the design but rather the implementation. Ultimately there was not enough up to date information about the available links' jitter, so that the path selection algorithm could choose the right one.

4.6 Improvements

There are several possible improvements that immediately spring to mind. These could be applied both to the actual implementation as well as to the high level design.

Measuring and Predicting Path Characteristics

Firstly, the path estimation could be adjusted to not just report previous statistics but also attempt to infer what the path might look like in the near future, for example as a path begins to experience increased latency a predictive algorithm/approach might be able to preemptively move flows off of that path, before their latency requirements are violated. This approach could be based on machine learning or AI, or it could use analytical methods.

Furthermore, it would be prudent to use an exponentially weighted averaging function when calculating a path's current characteristics, based on previous measurements. This would strike a nice balance between re-actively adjusting the predictions based on previous measurements and keeping in mind a path's performance across its entire history.

Also, as the authors in [19] did, this thesis' approach could benefit from adaptive windows of reporting. Instead of collecting statistics at a constant periodic rate, the rate can be adjusted to that period which performs best. This way there is not additional overhead with overly frequent statistical updates, as well as avoiding the reverse situation, where the reporting is too infrequent for a rapidly changing path.

Lastly, the probe packets which are sent out periodically to measure each path could be adjusted to periodically be sent in large bursts, this would serve to provide a more accurate idea of the jitter and the packet loss, since these cannot be measured purely by sending a small number of probe packets. For accurate jitter and reliability measurements more packets are required. This may incur larger overhead on the link due to the bursts of extra packets, however one idea to counteract this would be to adjust the probing so that it is only performed on those links with insufficient traffic currently running through them. The links which are currently backhauling large amounts of traffic will be able to make their own measurements based purely on the backhaul traffic, and there may be no need at all for probe packets. Conversely, the underutilized links will not be able to make good estimations, and sending bursts of probe packets will not interfere with the actual traffic. This would present a solution to one of the biggest issues with the approach presented in this thesis.

Adjustments for Jitter

Another potential improvement would be to pass a "Time of Execution" field in the GTP header of packets of jitter-sensitive applications. This allows the receiving WAN connector to store packets if they have arrived too early, and, conversely, the packet ordering function may use this field to determine that it is more important to forward the current packet now, than to wait for a missing packet. This store and forward approach, with the time of execution field, is one of the solutions mentioned in the Deterministic Networking specification and would be a very sensible addition to this solution.

Additionally, specifically for the Lower Earth Orbit (LEO) satellite case, the path selection could potentially be adjusted to account for the periodic increases in latency as the current satellite leaves the range of the ground station, and the next one comes into range. For example during this phase it might make sense to temporarily forward packets on a different path, or to enable the store and forward mechanism, before disabling it once the satellite handover has passed.

Forward Error Correction

One possible improvement is the addition of Forward Error Correction (FEC). While this would be difficult to integrate into the equation to select paths, FEC could be used to increase the resilience of consistently lossy links, which is a big benefit for links which commonly exhibit this characteristic, such as wireless links. Especially LEO satellite links tend to present around 1% packet loss [12], which makes the use of FEC much more reasonable since one would only need to add a 2 to 5% overhead, using a fountain code scheme, and thus wind up with a much more reliable link. However as was observed in the evaluation of the WAN Connector's throughput when using the Packet Ordering and Elimination Function, additional computation when receiving packets can lead to reduced overhead. To this extent, the addition of FEC would require efficient decoders and encoders on the sending and receiving side, so that the packet processing speed is not slowed down.

Traffic Prioritization

Perhaps the greatest theoretical flaw in the presented approach is that the WAN Connector does not perform packet prioritization, even though the 5G specifications explicitly mention that flows may have different priorities and should be prioritized accordingly. This is a significant flaw, and very difficult to account for. The most reasonable approach to integrating this is to add it to the traffic shaper. In the experiments performed here, the CAKE traffic

shaper was used. It's implementation in the linux kernel does allow for up to 8 different priority "tins" for traffic, however that is not enough for the 5G requirements which allow for a great deal more priority levels than just 8. One option, therefore, would be to extend the CAKE implementation to allow a greater number of priority tins. Alternatively, the a different traffic shaper could be chosen, or developed, which incorporates prioritization as well as fairness and active queue management at the level of individual IP flows.

Load Balancing

In instances where more than one path is viable, the algorithm used for path selection does not account for the current load of the paths, nor does it try to perform load balancing when selecting paths. This is because the algorithm attempts to find the minimum weight path. Even in cases where two paths have the same properties and the same weighting, the algorithm does not attempt to shuffle or rotate between them. All flows with similar requirements will wind up on the same path. This can lead to under utilization of other paths and congestion problems on the "primary" path. This could be partially addressed by integrating a degree of random selection for equal paths. However a potentially more interesting approach would be to dynamically update the weights of different paths, based on their current load. This would keep randomness out of the algorithm, for reproducible results, and act as a form of load balancing. However great care would have to be taken with respect to how the weights are adjusted, so that paths are still appropriately weighted relative to each other. i.e. a path which is deemed twice as expensive as another path should not have more than half as much traffic as the other path, for flows which could go be forwarded on either path.

Chapter 5

Conclusion and Outlook

5.1 Summary

In review, this thesis has addressed the usage of multiple backhaul paths in a campus 5G setting to achieve determinism. A review of the topic and relevant literature was done in the background chapter, then an approach was developed based on that information and based on the requirements of the problem. Finally this approach was evaluated in a testbed which emulated multiple outgoing links, to see how it performed. The implementation provided was able to achieve some of the goals set for latency and packet loss, but struggled with jitter reduction, as well as utilizing all of the available bandwidth when it is performing packet replication to improve reliability. In the discussion of these results, their potential sources and improvements were addressed.

Placed in a wider context, this thesis presents a foray into an area which will become more and more relevant as 5G campus deployments increase and mature, and as 5G applications which require greater degrees of determinism become more commonplace. On the whole, greater degrees of bespoke packet processing are required for traffic on the internet, as application's evolve and their requirements become more strict. Especially in geographically distributed

settings with multiple backhaul paths, the ability to intelligently select among these paths will be crucial to enabling these types of applications.

5.2 Implications and Further Areas of Research

The results obtained here imply the presented approach is worthy of further investigation, however only in conjunction with some of the presented improvements. The implementation provided by this thesis cannot achieve determinism without certain adjustments, but nor can it be concluded that it will be definitively unable to do so.

Since the results were only verified in a testbed setup it would make sense to now test the WAN Connector in a real campus 5G deployment where there actually are multiple outgoing paths. Furthermore many of the improvements suggested in the previous section would all be worthy of implementing and evaluating in similar experiments.

Research should also be conducted on evaluating the performance of specific applications performance. For example, the quality of VoIP calls doesn't depend just on latency and jitter [36], but rather how they interact together, and they have their own suites of evaluation criteria. Another specific application to consider has to be interactive video. Mission critical as well as control systems could also be considered in the test suites for evaluation.

It would be interesting to see what benefit AI and machine learning approaches may bring to this problem since they can act more dynamically, and perhaps learn the characteristics of a given link over time. Perhaps they can discover what characteristics the link exhibits right before total failure and thus perform pre-emptive path switching.

Appendix A

Patch for the BPF Flow Dissector to Peek Inside GTP

As was mentioned in the "Approach" chapter, in a true 5G Campus deployment the messages between either the RAN nodes or the UPF and the WAN Connector will be carrying GTP packets. In order for the tc-cake traffic shaper in the linux kernel to be able to detect the packets being tunneled, a custom BPF flow dissector must be used ¹.

The linux kernel provides a very simple example of a set of dissectors which remove the tunnel headers for various protocols. Unfortunately the GTP protocol is not included in this set, but it is a very simple patch to add the required functionality to the kernel.

The listing [A.1](#) shows a patch for the example file given in the linux kernel ². This patch constitutes an important part of this thesis' contribution because without it, a realistic deployment of the proposed system would not be possible, because the traffic shaper would not function.

¹https://www.kernel.org/doc/html/v5.1/networking/bpf_flow_dissector.html#overview

²https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/tree/tools/testing/selftests/bpf/progs/bpf_flow.c?h=

```

1
2 diff --git a/bpf_flow.c b/bpf_flow.c
3 index ed437e1..70a08a8 100644
4 --- a/bpf_flow.c
5 +++ b/bpf_flow.c
6 @@ -161,6 +161,8 @@ int _dissect(struct __sk_buff *skb)
7         return parse_eth_proto(skb, keys->n_proto);
8     }
9
10 #define GTP_ENCAP_OFFSET 24
11 +
12 /* Parses on IPPROTO_* */
13 static __always_inline int parse_ip_proto(struct __sk_buff *skb
14     , __u8 proto)
15 {
16 @@ -171,6 +173,10 @@ static __always_inline int parse_ip_proto(
17     struct __sk_buff *skb, __u8 proto)
18     struct ethhdr *eth, _eth;
19     struct tcphdr *tcp, _tcp;
20     struct udphdr *udp, _udp;
21 +    struct udphdr *encap_udp, _encap_udp;
22 +    struct iphdr *iph, _iph;
23 +    __u16 srcport;
24 +    __u16 destport;
25
26     switch (proto) {
27     case IPPROTO_ICMP:
28 @@ -242,7 +248,26 @@ static __always_inline int parse_ip_proto(
29     struct __sk_buff *skb, __u8 proto)
30         udp = bpf_flow_dissect_get_header(skb, sizeof(*
31         udp), &_udp);
32         if (!udp)
33             return export_flow_keys(keys, BPF_DROP);
34         srcport = bpf_ntohs(udp->source);
35         destport = bpf_ntohs(udp->dest);
36         if (srcport == 2152 && destport == 2152) {
37             keys->thoff += sizeof(*udp) +
38             GTP_ENCAP_OFFSET;
39             iph = bpf_flow_dissect_get_header(skb,
40             sizeof (*iph), &_iph);
41             if (iph) {
42                 keys->thoff += (iph->ihl << 2);
43                 encap_udp =
44                 bpf_flow_dissect_get_header(skb, sizeof(*encap_udp), &
45                 _encap_udp);
46                 if (encap_udp) {
47                     keys->thoff -= (iph->ihl
48                     << 2);

```

```

40 +                                     keys->is_encap = true;
41 +                                     return parse_eth_proto(
      skb, bpf_htons(ETH_P_IP));
42 +                                     } else {
43 +                                     bpf_printk("Couldnt
      dissect encapsulated udp header");
44 +                                     }
45 +                                     } else {
46 +                                     bpf_printk("Couldnt
      dissect encapsulated ip header at GTP offset (%d bytes)",
      keys->thoff);
47 +                                     }
48
49 +     }
50     keys->sport = udp->source;
51     keys->dport = udp->dest;
52     return export_flow_keys(keys, BPF_OK);

```

Listing A.1: Patch to remove GTP Headers for the BPF Flow Dissector Example in the Linux Kernel

Appendix B

JSON Formats Used in Control and Data-Plane Messages

This appendix contains excerpts of the JSON messages passed between the control plane and the data plane during operation. Listing [B.1](#) shows how the control plane sends updates to the data plane nodes. Listing [B.2](#) shows the format of the messages used for reporting statistics.

These are the only two message formats used between the Control Plane and the Data Plane. These messages are sent over SCTP, UTF-8 encoded. They are the only data in the SCTP messages' body. By using JSON there are three distinct advantages- firstly, no specific protocol must be developed for the messaging between control and data plane, secondly the messages are both human and machine readable which makes for easier debugging, and lastly because it is a hugely popular data format there are lots supporting libraries.

The messages which add new flows to the Data Plane contain the flow Description as well as its priority, and the Packet Ordering Function (POF) maximum delay, which was mentioned in the description of the POF algorithm in the approach.

In the reverse direction, when the Data Plane nodes report to the Control Plane, the statistics reported include firstly the paths description (source and destination IP address) as well as the number of packets and bytes received on that path, and the average jitter and latency experienced, both over the entire reporting period, and since the last message. Finally, the message are also sent with a list of the "sending" statistics. That is, the number of bytes and packets sent out on the WAN Connector's interfaces. Both the path and interface statistics are sent with three timestamps. One timestamp to describe the first packet sent or received during this reporting window, one to describe the last packet sent or received during that same reporting window, and a timestamp taken before the message itself was sent.

By including the "sending" statistics, the WAN Connector is able to inform it's own Control Plane if an interface is being actively used. On the receiving side, if an interface has not received packets from the respective source in too long a time, then the Control Plane can infer complete link failure. This is crucial for the ability of the WAN Connector to switch all user traffic off of dead links. Once a link works again, the WAN Connector can detect this because the periodic probe packets will be received.

In order to save space this message has been edited for display so that it only contains messages about two of the available paths instead of all the paths.

```

1
2 {
3   "flowsToAdd": [
4     {
5       "flowID": 1,
6       "priority": 1,
7       "flowDescription": {
8         "src": "10.0.0.31",
9         "srcSubnet": 24,
10        "srcPortStart": 0,
11        "srcPortEnd": 0,
12        "dst": "10.0.3.31",
13        "dstSubnet": 24,
14        "dstPortStart": 0,
15        "dstPortEnd": 0,
16        "protocol": 1,
17        "direction": "in"
18      },
19      "POFmaxDelayMicros": 0,
20      "nForwarding": 1,
21      "decisions": [
22        {
23          "interfaceName": "link3",
24          "remoteEndpoint": {
25            "ip": "10.1.22.21",
26            "port": 2152
27          }
28        }
29      ]
30    },
31    {
32      "flowID": 2,
33      "priority": 1,
34      "flowDescription": {
35        "src": "10.0.0.31",
36        "srcSubnet": 24,
37        "srcPortStart": 0,
38        "srcPortEnd": 0,
39        "dst": "10.0.3.31",
40        "dstSubnet": 24,
41        "dstPortStart": 0,
42        "dstPortEnd": 0,
43        "protocol": 17,
44        "direction": "in"
45      },

```

```

46     "POFmaxDelayMicros": 0,
47     "nForwarding": 1,
48     "decisions": [
49         {
50             "interfaceName": "link3",
51             "remoteEndpoint": {
52                 "ip": "10.1.22.21",
53                 "port": 2152
54             }
55         }
56     ],
57 },
58 {
59     "flowID": 3,
60     "priority": 1,
61     "flowDescription": {
62         "src": "10.0.0.31",
63         "srcSubnet": 24,
64         "srcPortStart": 0,
65         "srcPortEnd": 0,
66         "dst": "10.0.3.31",
67         "dstSubnet": 24,
68         "dstPortStart": 0,
69         "dstPortEnd": 0,
70         "protocol": 6,
71         "direction": "in"
72     },
73     "POFmaxDelayMicros": 0,
74     "nForwarding": 1,
75     "decisions": [
76         {
77             "interfaceName": "link3",
78             "remoteEndpoint": {
79                 "ip": "10.1.22.21",
80                 "port": 2152
81             }
82         }
83     ]
84 }
85 ]
86 }

```

Listing B.1: Format of Flow Forwarding Decisions from Control Plane


```

1
2 {
3   "pathStats": [{
4     "src": "10.1.33.31",
5     "dst": "10.0.33.32",
6     "nbytes": 6272,
7     "npackets": 224,
8     "ndropped": 0,
9     "avgLatencyMicros": 30655,
10    "avgJitterMicros": 1835,
11    "bytesSince": 56,
12    "packetsSince": 2,
13    "droppedSince": 0,
14    "latencySinceMicros": 29276,
15    "jitterSinceMicros": 1431,
16    "earliestRecvMillis": 1703448057713,
17    "lastRecvMillis": 1703448057716,
18    "heartbeatTimestampMillis": 1703448058929
19  }, {
20    "src": "10.1.44.41",
21    "dst": "10.0.44.42",
22    "nbytes": 6272,
23    "npackets": 224,
24    "ndropped": 0,
25    "avgLatencyMicros": 28648,
26    "avgJitterMicros": 4182,
27    "bytesSince": 56,
28    "packetsSince": 2,
29    "droppedSince": 0,
30    "latencySinceMicros": 29540,
31    "jitterSinceMicros": 5351,
32    "earliestRecvMillis": 1703448057711,
33    "lastRecvMillis": 1703448057718,
34    "heartbeatTimestampMillis": 1703448058929
35  }],
36  "interfaceStats": [{
37    "name": "link2",
38    "src": "10.0.33.32",
39    "nbytes": 162768024,
40    "npackets": 109020,
41    "bytesSince": 56,
42    "packetsSince": 2,
43    "earliestSentMillis": 1703448055928,
44    "lastSentMillis": 1703448055928,
45    "heartbeatTimestampMillis": 1703448058929

```

```

46     }, {
47         "name": "link4",
48         "src": "10.0.44.42",
49         "nbytes": 174418872,
50         "npackets": 116808,
51         "bytesSince": 56,
52         "packetsSince": 2,
53         "earliestSentMillis": 1703448055928,
54         "lastSentMillis": 1703448055928,
55         "heartbeatTimestampMillis": 1703448058929
56     }]
57 }

```

Listing B.2: Format of Statistics from Data Plane

Appendix C

Example of Custom GTP Header Format

This appendix contains a screenshot of a packet captured during normal operation of the WAN Connector (figure C.1), being displayed using wireshark ¹. The packet dissection shows the GTP format being used, along with the additional custom header, which is 8 bytes long. The custom header consists of 2 bytes of regular overhead from the GTP protocol as well as a 4 byte timestamp and 2 spare bytes. This packet was captured during the jitter experiment.

¹<https://www.wireshark.org/>

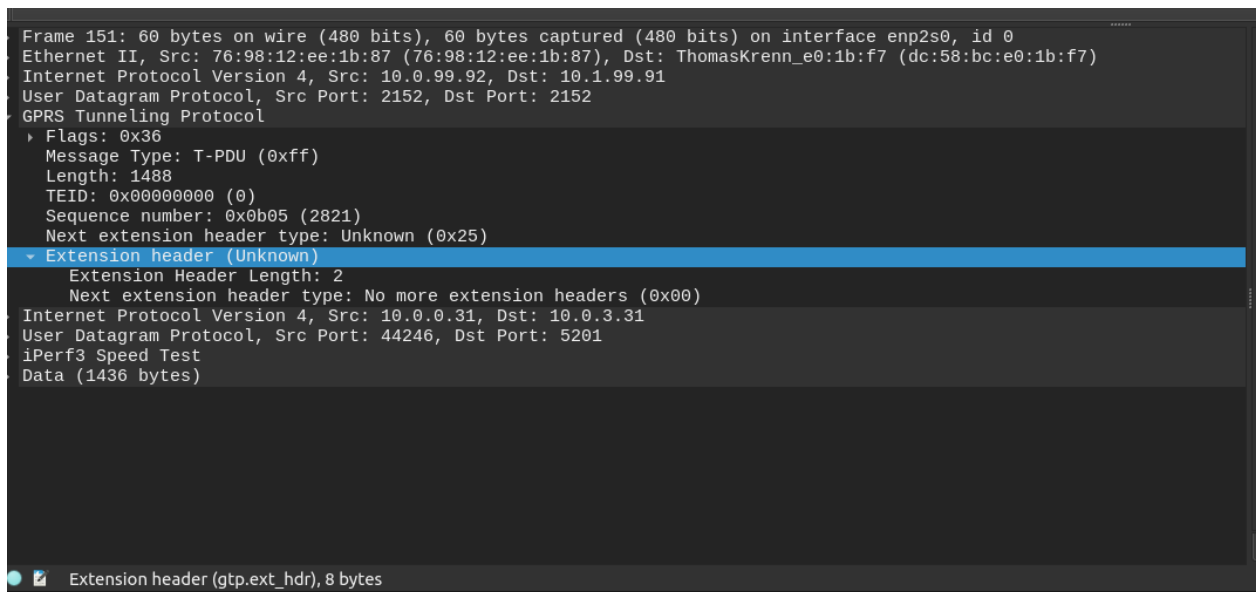


Figure C.1: Screenshot of Wireshark Dissection

Bibliography

- [1] 3GPP. *Evolved Universal Terrestrial Radio Access (E-UTRA) and Evolved Universal Terrestrial Radio Access Network (E-UTRAN); Overall description; Stage 2*. Technical Specification (TS) 36.300. Version 17.4.0. 3rd Generation Partnership Project (3GPP), Apr. 2023.
- [2] 3GPP. *General Packet Radio Service (GPRS); GPRS Tunnelling Protocol (GTP) across the Gn and Gp interface*. Technical Specification (TS) 29.060. Version 17.4.0. 3rd Generation Partnership Project (3GPP), Oct. 2022.
- [3] 3GPP. *Study on architecture aspects for using satellite access in 5G*. Technical Report (TR) 23.737. Version 17.2.0. 3rd Generation Partnership Project (3GPP), Mar. 2021.
- [4] 3GPP. *System architecture for the 5G System*. Technical Specification (TS) 23.501. Version 17.8.0. 3rd Generation Partnership Project (3GPP), Apr. 2023.
- [5] Rudolf Ahlswede et al. “Network information flow”. In: *IEEE Transactions on information theory* 46.4 (2000), pp. 1204–1216.
- [6] Aditya Akella et al. “A measurement-based analysis of multihoming”. In: *Proceedings of the 2003 conference on Applications, technologies, architectures, and protocols for computer communications*. 2003, pp. 353–364.
- [7] Aditya Akella et al. “On the performance benefits of multihoming route control”. In: *IEEE/ACM Transactions on Networking* 16.1 (2008), pp. 91–104.

- [8] Mark Allman. “Comments on bufferbloat”. In: *ACM SIGCOMM Computer Communication Review* 43.1 (2012), pp. 30–37.
- [9] Hind Alwan and Anjali Agarwal. “Multi-objective reliable multipath routing for wireless sensor networks”. In: *2010 IEEE Globecom Workshops*. IEEE. 2010, pp. 1227–1231.
- [10] Austin Appleby. “MurmurHash3, 2012”. In: *URL: <https://github.com/aappleby/smhasher/blob/master/cpp>* (2012).
- [11] Anat Bremner-Barr et al. “Predicting and bypassing end-to-end Internet service degradations”. In: *Proceedings of the 2nd ACM SIGCOMM Workshop on Internet Measurement*. 2002, pp. 307–320.
- [12] Joerg Deutschmann, Kai-Steffen Hielscher, and Reinhard German. “Broadband internet access via satellite: Performance measurements with different operators and applications”. In: *Broadband Coverage in Germany; 16th ITG-Symposium*. VDE. 2022, pp. 1–7.
- [13] Norman Finn and Pascal Thubert. *Deterministic Networking Problem Statement*. RFC 8557. May 2019. DOI: [10.17487/RFC8557](https://doi.org/10.17487/RFC8557). URL: <https://www.rfc-editor.org/info/rfc8557>.
- [14] Norman Finn et al. *Deterministic Networking Architecture*. RFC 8655. Oct. 2019. DOI: [10.17487/RFC8655](https://doi.org/10.17487/RFC8655). URL: <https://www.rfc-editor.org/info/rfc8655>.
- [15] Igor Ganichev et al. “YAMR: Yet another multipath routing protocol”. In: *ACM SIGCOMM Computer Communication Review* 40.5 (2010), pp. 13–19.
- [16] Jim Gettys. “Bufferbloat: Dark buffers in the internet”. In: *IEEE Internet Computing* 15.3 (2011), pp. 96–96.
- [17] David K Goldenberg et al. “Optimizing cost and performance for multihoming”. In: *ACM SIGCOMM Computer Communication Review* 34.4 (2004), pp. 79–92.
- [18] Ethan Grossman. *Deterministic Networking Use Cases*. RFC 8578. May 2019. DOI: [10.17487/RFC8578](https://doi.org/10.17487/RFC8578). URL: <https://www.rfc-editor.org/info/rfc8578>.

- [19] Ahsan Habib and John Chuang. “Improving application QoS with residential multihoming”. In: *Computer Networks* 51.12 (2007), pp. 3323–3337.
- [20] Toke Høiland-Jørgensen, Dave Täht, and Jonathan Morton. “Piece of CAKE: a comprehensive queue management solution for home gateways”. In: *2018 IEEE International Symposium on Local and Metropolitan Area Networks (LANMAN)*. IEEE. 2018, pp. 37–42.
- [21] Xiaoxia Huang and Yuguang Fang. “Multiconstrained QoS multipath routing in wireless sensor networks”. In: *Wireless Networks* 14.4 (2008), pp. 465–478.
- [22] Mona Jaber et al. “5G backhaul challenges and emerging research directions: A survey”. In: *IEEE access* 4 (2016), pp. 1743–1766.
- [23] Mark W Krentel. “The complexity of optimization problems”. In: *Proceedings of the eighteenth annual ACM symposium on Theory of computing*. 1986, pp. 69–76.
- [24] Michael Luby et al. *Raptor forward error correction scheme for object delivery*. Tech. rep. 2007.
- [25] Sami Ma et al. “Network characteristics of LEO satellite constellations: A Starlink-based measurement from end users”. In: *IEEE INFOCOM 2023-IEEE Conference on Computer Communications*. IEEE. 2023, pp. 1–10.
- [26] Rong Pan, Balaji Prabhakar, and Konstantinos Psounis. “CHOKe-a stateless active queue management scheme for approximating fair bandwidth allocation”. In: *Proceedings IEEE INFOCOM 2000. Conference on Computer Communications. Nineteenth Annual Joint Conference of the IEEE Computer and Communications Societies (Cat. No. 00CH37064)*. Vol. 2. IEEE. 2000, pp. 942–951.
- [27] Charles Perkins, Elizabeth Belding-Royer, and Samir Das. *RFC3561: Ad hoc on-demand distance vector (AODV) routing*. 2003.
- [28] Jonathan Prados-Garzon and Tarik Taleb. “Asynchronous time-sensitive networking for 5G backhauling”. In: *IEEE Network* 35.2 (2021), pp. 144–151.

- [29] Jonathan Prados-Garzon et al. “5G non-public networks: Standardization, architectures and challenges”. In: *IEEE Access* 9 (2021), pp. 153893–153908.
- [30] Yakov Rekhter, Tony Li, and Susan Hares. *A border gateway protocol 4 (BGP-4)*. Tech. rep. 2006.
- [31] Justus Rischke et al. “5G campus networks: A first measurement study”. In: *IEEE Access* 9 (2021), pp. 121786–121803.
- [32] Salman Saadat, Da Chen, and Tao Jiang. “Multipath multihop mmWave backhaul in ultra-dense small-cell network”. In: *Digital Communications and Networks* 4.2 (2018), pp. 111–117.
- [33] Kari Seppänen et al. “Multipath routing for mmWave WMN backhaul”. In: *2016 IEEE International Conference on Communications Workshops (ICC)*. IEEE. 2016, pp. 246–253.
- [34] Ion Stoica, Hui Zhang, and TS Eugene Ng. “A hierarchical fair service curve algorithm for link-sharing, real-time and priority services”. In: *ACM SIGCOMM Computer Communication Review* 27.4 (1997), pp. 249–262.
- [35] Shu Tao et al. “Exploring the performance benefits of end-to-end path switching”. In: *Proceedings of the joint international conference on Measurement and modeling of computer systems*. 2004, pp. 418–419.
- [36] Shu Tao et al. “Improving VoIP quality through path switching”. In: *Proceedings IEEE 24th Annual Joint Conference of the IEEE Computer and Communications Societies*. Vol. 4. IEEE. 2005, pp. 2268–2278.
- [37] Alok S Tongaonkar. “Fast pattern-matching techniques for packet filtering”. PhD thesis. Citeseer, 2004.
- [38] Balazs Varga et al. *Deterministic Networking (DetNet): Packet Ordering Function*. Internet-Draft draft-ietf-detnet-pof-08. Work in Progress. Internet Engineering Task Force, Dec. 2023. 12 pp. URL: <https://datatracker.ietf.org/doc/draft-ietf-detnet-pof/08/>.

- [39] Jaspreet Singh Walia, Heikki Hämmäinen, and Marja Matinmikko. “5G Micro-operators for the future campus: A techno-economic study”. In: *2017 internet of things business models, users, and networks*. IEEE. 2017, pp. 1–8.