

TCP PEP

TCP Performance Enhancing Proxy to Support
Non-interactive Applications

Joe Bayer

Informatics: Programming and System Architecture
60 ECTS study points

Department of Informatics
Faculty of Mathematics and Natural Sciences

Joe Bayer

TCP PEP

TCP Performance Enhancing Proxy to
Support Non-interactive Applications

Supervisors:
Michael Welzl
Kristjon Ciko

Contents

1	Intro	3
2	Essay	4
2.1	TCP/IP	4
2.1.1	Congestion Control	5
2.1.2	3 Way handshake (0 RTT)	6
2.1.3	TCP Options and Fast Open	7
2.2	Future of wireless communication.	8
2.2.1	5G Millimetre Wave	8
2.2.2	Buffering	9
2.2.3	Non-Interactive Applications	10
2.3	Proxy	10
2.3.1	PEP	11
2.3.2	PEP for wireless communication	11
2.3.3	Transparent vs Non-Transparent	11
2.4	Linux	12
2.4.1	Kernel Modules	12
3	Design	13
3.1	Connection Splitting	13
3.2	Connection in 0 round trips	13
3.3	Transparent vs Non-Transparent	13
3.4	Kernel Module Vs. Userspace Application	13
3.5	Summary	13
4	Implementation	15
4.1	Kernel Module	15
4.2	TLV Library	16
4.3	PEP	16
4.4	PEP Connect	16
5	Evaluation	17
5.1	Initial configuration	17
5.1.1	Traffic Control Options	17

5.1.2	Scheduling Algorithms	17
5.1.3	Interactive vs PEP	18
5.1.4	PEP vs E2E tests	18
5.1.5	10x10 tests	19
5.1.6	Spike	19
6	Conclusion	20
6.1	Future Work	20

Chapter 1

Intro

Chapter 2

Essay

In this chapter we will present some of the required background knowledge to understand the concepts presented in this paper. Focusing on topics that are outside the common understanding of network programming, especially details of certain congestion controllers and network protocols will be discussed. The rest of the thesis will assume the following topics are known to the reader.

2.1 TCP/IP

Perhaps the most well known internet transport protocol is the Transmission Control Protocol (TCP). It is known for providing reliable and in-order delivery of packets using acknowledgments and re-transmissions [5]. It was first introduced in 1974, but is still the most used internet protocol. However, as the demands of the internet have changed, TCP has not. Though TCP has been updated with minor extensions over the years, such as increased initial windows or new options, the core ideas have stayed the same [2].

Concepts as the end-to-end argument still play a vital role in how TCP is used in the modern internet. TCP is suffering under the illusion that all logic should be placed on the endpoints as the end to end argument denotes. Even if it spans multiple different domains with varying topologies and demands, especially between wired and wireless domains.

- **Wireless Domain:** A wireless communication domain refers to the transmission of data over a wireless medium without the use of physical connections such as wires or cables between devices. This domain covers a variety of technologies, including 3G, 4G, and 5G for mobile communication, Bluetooth and Wi-Fi for close-range communication, and satellite communication for worldwide communication.
- **Wired Domain:** Unlike a wireless domain, a wired domain provides a

steady and reliable bandwidth with low error rates and high throughput. The use of Ethernet and Fiber are typical for wired networks, they enable the transmission of large amount of data over long distances and low signal noise.

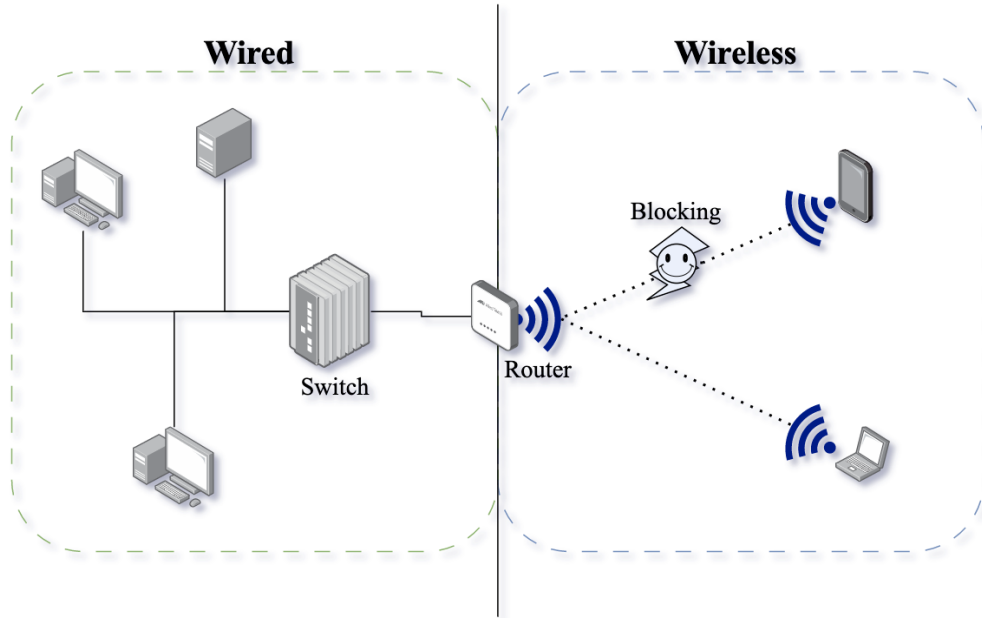


Figure 2.1: Example of network domains

Each domain has different requirements that a single TCP connection cannot provide. Fig. 2.1 shows the two domains and their characteristic differences. Usually, wireless domains experience a lot of changes in connectivity and bandwidths, while the wired domain usually is considered stable. This creates problems for the "modern" TCP which, because of the end to end argument, normally spans multiple domains. Especially congestion control has problems adapting to high fluctuating bandwidth across long distances and multiple domains.

2.1.1 Congestion Control

Congestion occurs in the internet when a network's resources, such as routers, are overloaded to the point that they diminish quality of the network [11]. Packet loss and high delays are common issues associated to high congestion in the network. To solve the problem of congestion, a distributed algorithm is used: Congestion Control. The main goal of congestion control is to maintain a stable network, while still utilizing the available bandwidth shared among all flows. This is achieved by additively increasing the sending rate, and multiplicatively reducing the sending rate when detecting congestion [10].

Congestion can be detected by packet loss, changes in delay, but also by explicit notifications.

Over time different variations of congestion controller have emerged. Although their goal is the same, reduce congestion in the network, their approaches vary.

- **TCP Reno:** Reno embodies the traditional approach to congestion control. Slowly increasing the sending rate while the network is stable and drastically reducing it on packet loss. TCP Reno was designed for unstable and dynamic networks, where the rapid response rate is crucial to prevent network overloading. However, the slow start rate and aggressive reduction of the sending rate make it sub optimal for more stable networks, where packet loss is less frequent and predictable. Consequently, TCP Reno's reliance on packet loss may lead to unnecessary rate reductions and decreased network throughput.
- **New Vegas:** New Vegas is similar to TCP Reno in most aspects, the main difference is the use of delay to detect congestion instead of packet loss. This makes New Vegas able to react faster to congestion, however it also introduces some interesting side effects. If New Vegas competes with TCP Reno flows, it will start reducing its sender rate before TCP Reno does, this leads to New Vegas losing out on possible bandwidth.
- **Cubic:** Cubic improves on the idea of TCP Reno by using a cubic function to adjust its (congestion window) sending rate in order to achieve higher throughput in a fast manner. Cubic is very efficient in highspeed networks and known for handling large data transfer over long distances. However, Cubic is not as reliable and robust as more traditional congestion controllers like TCP Reno.

In summary, the main differences between TCP Reno, New Vegas and Cubic are their approach to congestion control, their performance in different types of networks, and their trade-off between efficiency and reliability.

2.1.2 3 Way handshake (0 RTT)

For TCP to establish a connection it uses a three-way handshake. Initially, it transmits a synchronization (SYN) packet to the desired endpoint. The endpoint responds with an acknowledgement and a synchronization packet of its own (SYN/ACK). Finally, the client responds with a acknowledgment (ACK). At this point both endpoints have confirmed that they are ready for further communication. For any connection to be established this handshake has to be done. For short flows that terminate in just a few round trips

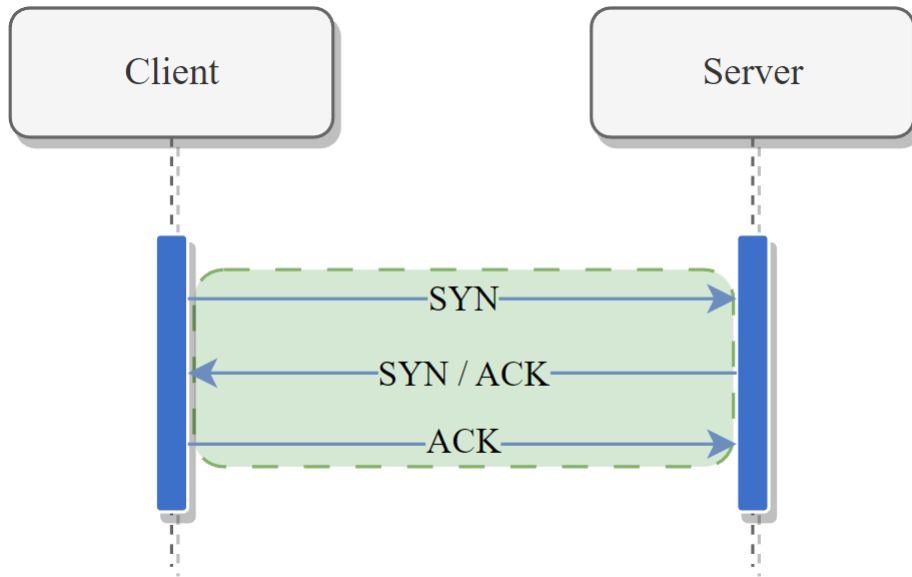


Figure 2.2: The TCP handshake procedure

the initial TCP handshake can be a bottleneck, which is made worse if the connection is using a proxy and has to exchange additional information.

2.1.3 TCP Options and Fast Open

A TCP connection can be configured with optional header extensions called TCP Options [3]. These options change the default behaviour of TCP or add new features. One such feature is TCP Fast Open, which allows data to be added to the initial synchronization packet. A typical use case could be adding a HTTP GET request, thereby saving an entire round trip. In general flows that terminate in a few round trips greatly benefit from this feature. The reason being, the bottleneck in such connections often lies within the initial TCP handshake. Therefore, by removing the extra round trip required to send the first data packet, a significant amount of time can be saved.

TCP Fast Open also has other benefits, such as establishing connections to proxies [2]. When you are trying to establish a connection through a proxy, you get the added delay of a second round trip sending the desired endpoint. This can be avoided by using TCP Fast Open to send the desired endpoint in the first synchronization packet to the proxy. "SYN forwarding" enables the users to establish a proxy connection without any added delays, however it does depend on the user's application to use TCP Fast Open.

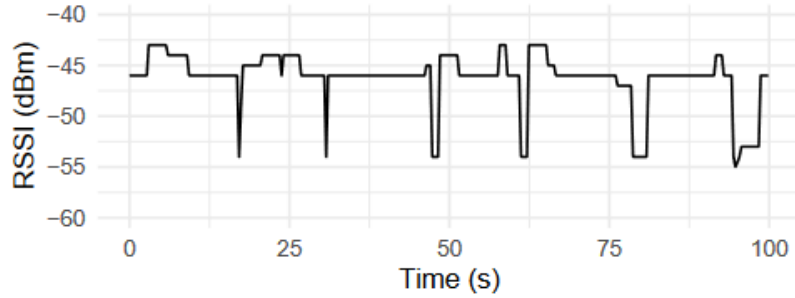


Figure 2.3: 5G bandwidth fluctuations from humans

2.2 Future of wireless communication.

The future of wireless communication has seen a lot of improvements such as highly increased bandwidth achieved through advanced technologies like 5G and beyond. Millimetre frequency bands have opened up new possibilities for wireless communication. These higher frequency bands offer greater capacity and can accommodate more devices, however high frequencies come with a set of new challenges such as highly fluctuating bandwidths. This fluctuation can be influenced by various factors such as signal interference, obstacles in the signal path, and environmental conditions.

2.2.1 5G Millimetre Wave

The emergence of 5G Millimeter wave communications has opened the doors for low latency networks with multiple gigabit bandwidth. This is achieved by using higher millimetre wave (mmWave) frequencies in the range of 30GHz to 300GHz, which has a lot of benefits [1]. A wider spectrum of frequencies to choose from and higher data transfer rates are just some of the many benefits mmWave provides. But along side the benefits, mmWave has also introduced a lot of new challenges.

A big problem with millimetre wave communication is signal path blocking also called "Line of sight blocking" [8]. It's caused by the use of Beam-forming to increase the bandwidth and range of millimeter wave signals. Beam-forming focuses the signal in a certain direction making any blocking of the signal path devastating for the bandwidth. Even the human body can create enough blockage to drastically reduce the bandwidth. This causes huge fluctuations in the bandwidth whenever the signal is blocked.

Fluctuating bandwidths lead to unstable TCP connections with a worst case of losing packets. Current TCP congestion controllers such as CU-BIC, New Reno or New Vegas struggle when reacting to sudden fluctuating

changes. They are simply not able to utilize the high bandwidth when it is available. Simply increasing the aggressiveness of a congestion controller is not an option either as it would disrupt the internet and not be TCP-friendly. A possible solution could be to buffer packets at the 5G base stations, having the data ready for when the bandwidth is high. This however creates a new problem, bufferbloat.

2.2.2 Buffering

Buffer bloat

The buffer bloat problem occurs when the systems between the endpoints buffer so many packets that the latency drastically increases and the reliability of the network as a whole goes down. The increased latency is detrimental for interactive (latency sensitive) applications. Generally it's preferred to drop packets and keep buffers small to avoid buffering time sensitive packets such as synchronization packets. Although this works in most cases, it's far from an optimal solution.

The increased bandwidth and low latency promises of new technology such as 5G has put a lot of pressure on the efficient forwarding of packets. Small buffers are therefore the standard, But at the same time, fluctuating bandwidth has shown the potential need to buffer packets for non-interactive traffic. Most focus has been on accommodating for latency sensitive applications like virtual reality or remote surgery to name a few.

This thesis will explore non-interactive applications where latency is not that critical and more buffering is acceptable and most likely desirable. By splitting traffic into interactive and non-interactive we can improve the performance of both. By having very small buffers for interactive applications we avoid bufferbloat problems, while utilizing the benefits of big buffers for non-interactive applications.

Packet Scheduling

A method of reducing the effects of bufferbloat is packet scheduling. A system should not send more packets than the weakest link can handle, this idea is built into TCP in the form of congestion control. However, when buffers grow to the point of causing bufferbloat, TCP's congestion control algorithms are unable to confidently determine a sending rate. Packet scheduling can solve this problem as it usually controls the size of the buffers. It makes sure queues can grow when needed, but keep the overall state of the buffers low. Packet scheduling has a lot more to offer than simple queue management,

this will be explored later.

Proposed packet scheduling algorithms:

- **FQ CoDel:** The Flow Queue Controlled Delay algorithm, FQ CoDel for short, was partly developed to deal with the bufferbloat problem. Its main goal is to reduce the impact of head-of-line blocking and give a fair share of bandwidth by mixing packets from multiple flows [6]. Internally FQ CoDel uses a FIFO queue, classifying packets into different flows to provide a fair share of bandwidth.
- **HTB:** Hierarchical Token Bucket is a queuing discipline based on assigning different classes a certain amount of bandwidth and sending rate. Because of its extensive bandwidth and delay management it's a good option for testing, especially in a virtual environment.

2.2.3 Non-Interactive Applications

Non-Interactive applications such as web traffic, file transfers and video streaming can benefit from larger buffering, especially with fluctuating bandwidths. This is because if we are able to buffer the packets closer to their final destination, we have them ready to be sent when the bandwidth changes. By buffering them we can decrease delay times and achieve faster total completion times for non-interactive traffic.(need citation or prove it myself?). At the same time, interactive applications will not suffer under large queue delays that occur under normal buffering.

2.3 Proxy

Proxy servers play a big role in the modern internet, delivering benefits such as anonymity and increased performance [9]. A common use case for a proxy is caching by keeping a copy of popular resources such as a websites. This reduces the latency of accessing the resource as long as the proxy is closer to the user than the original copy. Locality plays an important role in the total latency as any transmission will always be limited by the speed of light.

A proxy can also be used for privacy similar to a Virtual Private Network (VPN). By redirecting network traffic through a proxy, the origin of the traffic appears to be the proxy server rather than the actual end-user. Hypertext Transfer Protocol (HTTP), a popular internet protocol used for accessing websites, has this functionality built in using HTTP tunnels and a special CONNECT method in its header.

```
CONNECT mn.uio.no/:22 HTTP/1.1
```

2.3.1 PEP

A performance enhancing proxy (PEP) is a connection splitting proxy designed to increase performance of applications using it. The idea behind the PEP is putting more logic such as connection management, buffering, caching inside the network. As the name suggests, a PEP is designed to enhance the performance, but can also introduce new features to a network. An example of a new feature is the multipath support the TCP Transport converter gives. [2].

2.3.2 PEP for wireless communication

Performance enhancing proxies are already deployed and in use for a lot of wireless communication, especially satellites and radio access networks [7]. They have an inherent performance increase just by splitting the connection between the wireless and wired domains. These PEP's are therefore often installed at the base stations. However they are unable to distinguish between interactive or non-interactive traffic, meaning their buffers need to be low and still suffer from fluctuating bandwidth problems.

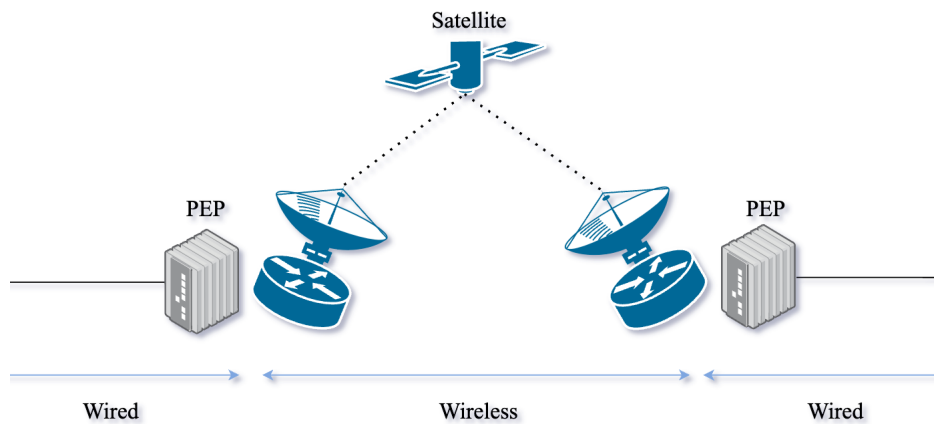


Figure 2.4: PEP installed to support Wireless traffic over satellite.

2.3.3 Transparent vs Non-Transparent

A big discussion regarding PEPs has been if they should be transparent or non-transparent. Transparent PEPs are not visible to the applications that use it. They silently split the connections and spoof the IP-address of both the client and server [4]. This is prone to cause unintended side effects,

such as certain TCP options not being forwarded and security concerns. Non-Transparent PEPs on the other hand are explicitly chosen by either the client or the server, and the sender is aware of the proxy splitting the original connection. This approach can be seen as more ethical and potentially remove some of the stigma associated with PEPs, this however requires modifications at the sender side utilize the PEP.

2.4 Linux

Linux is the most famous open source kernel freely available for anyone to use and modify. Because of the open source nature of Linux, there have been many various operating system implementation based on the Linux kernel. Ubuntu, Fedora or Manjaro are just some of the most famous Linux based operating systems out there. For developers, Linux is the perfect platform to experiment and test their new innovations. You are able to modify and recompile the kernel itself on the fly, and then test the solution on a live operating system. Linux supports most standards and is used by most major corporations such as Facebook, Amazon, Netflix and Google.

2.4.1 Kernel Modules

Thing that makes Linux truly extensible are Loadable Kernel Modules (LKM). Kernel modules are programmes that can be loaded at runtime into the kernel and run with kernel privileges. Running with kernel privileges has a lot of benefits such as having access to internal structures and kernel symbols. Most drivers in the Linux kernel are written as kernel modules as they need access to the system internals.

Congestion controllers and packet schedulers are also usually implemented as kernel modules. That is because Linux exposes a struct with function pointers that can be overwritten by a module. Making the kernel call the new functions instead. Because kernel modules run as part of the kernel they do not need to use system call to do basic I/O as using sockets. Removing the overhead of system calls makes the kernel modules run much faster than default userspace programs.

However, the using Linux kernel modules has the drawback that the program is bound to Linux. The modules will only work in the context of the Linux kernel as they depend on the internal functions, and that they are part of the kernel. Most other operating systems like MacOS will not allow user defined modules to run with kernel privileges. Additionally, any bugs or error in the kernel module will make the entire kernel panic, which usually requires a complete system restart to fix.

Chapter 3

Design

A good design for the PEP is crucial as it both needs to be robust, fast and reliable. Not all of the goals are equally achievable, and other aspects need to be considered such as cross platform compatibility. In this chapter we will explore the different ways of achieving our goals and compare them to each other.

3.1 Connection Splitting

A Performance enhancing proxy does not inherently need to be connection splitting. Based on where on the network stack it is implemented it could buffer and process packets without terminating the end to end connection. A PEP could simply listen to traffic and manipulate the traffic, or simply send ACKs early.¹

3.2 Connection in 0 round trips

3.3 Transparent vs Non-Transparent

3.4 Kernel Module Vs. Userspace Application

3.5 Summary

Table of design decisions based on different PEP implementations compared to ours. 0RTT, Transparent, TLVs, Special ACKS, connection splitting.

¹Picture of a PEP only sending ACKs

PEP List				
Implementation	0RTT	Connection Splitting	Special ACKs	Transparent
milliProxy	AF	AFG	004	x
PEPDNA	AX	ALA	248	x
SnoopTCP	AL	ALB	008	x
Our PEP	DZ	DZA	012	x
Transport Converter	AS	ASM	016	x
...	AD	AND	020	x
...	AO	AGO	024	x

Chapter 4

Implementation

4.1 Kernel Module

```
1 int pep_tcp_receive(struct socket *sock, u8* buffer, u32 size)
2 {
3     struct msghdr msg = {
4         .msg_flags = MSG_DONTWAIT,
5     };
6
7     struct kvec vec;
8     int rc = 0;
9
10    vec.iov_base = buffer;
11    vec.iov_len  = size;
12
13    printk(KERN_INFO "[PEP] kernel_recvmmsg: calling recvmmsg \n");
14    pep_tcp_receive_read_again:
15    rc = kernel_recvmmsg(sock, &msg, &vec, 1, vec.iov_len,
16        MSG_DONTWAIT);
17    if (rc > 0)
18    {
19        tlv_print(buffer);
20        printk(KERN_INFO "[PEP] kernel_recvmmsg: recvmmsg returned %d\n", rc);
21        return rc;
22    }
23    if(rc == -EAGAIN || rc == -ERESTARTSYS)
24    {
25        goto pep_tcp_receive_read_again;
26    }
27
28    printk(KERN_INFO "[PEP] kernel_recvmmsg: recvmmsg returned %d\n", rc);
29    return rc;
30 }
```

Listing 4.1: kernel_recvmmsg wrapper for receiving for TCP msgs

4.2 TLV Library

4.3 PEP

4.4 PEP Connect

Since PEP is explicitly addressed we need to change how we connect on the client side. The default approach is to use the **connect** function, our goal is to imitate this functions behavior, return code and parameters. By keeping our connect function as similar to the original we reduce the overhead of switching over to it. The only change to the original connect that we need is to specify type of traffic, interactive or non-interactive.

```
1 int socket, ret;
2 struct sockaddr_in s_in;
3 bzero((char *)&s_in, sizeof(s_in));
4 s_in.sin_family = AF_INET;
5 s_in.sin_addr.s_addr = inet_addr(IP);
6 s_in.sin_port = htons(PORT);
7
8 socket = socket(AF_INET, SOCK_STREAM, IPPROTO_TCP);
9
10 ret = connect(server, (struct sockaddr*) &s_in, sizeof(s_in));
11 ret = pep_connect(server, (struct sockaddr*) &s_in, sizeof(s_in),
    ), PEP_INTERACTIVE);
```

Listing 4.2: PEP connection

Chapter 5

Evaluation

5.1 Initial configuration

Sender -> Receiver -> Receiver

Data being pushed from Sender to receiver, file size, speed, delay.

5.1.1 Traffic Control Options

Linux has support for network interface configurations using the TC (traffic control) command. TC allows the configuration of packet scheduler, bandwidth, delay and jitter etc. These options combined with the fact that each network interface can have its own configuration, allows for very precise testing environments.

fq_codel does not allow the configuration of delay, this means we have to configure the delay on the path of the ACKS? We have to configure the delay on sender facing interface cards. ¹

5.1.2 Scheduling Algorithms

The choice of scheduling algorithm is very important for our test scenarios. As it will greatly affect the results of our tests, we will compare some against each other. The main goal is to highlight the effect and find the optimal algorithms for our PEP.

FIFO In our test case FIFO will have the effect of creating up a queue at our PEP. This is detrimental for the interactive UDP flow as it will be stuck in the queue, building up delay. This behavior can be explained by the fact that the non-interactive flow has a much higher sending rate than the interactive flow, thereby quickly filling the finite queue with non-interactive packets.

¹Picture of testbed setup

FQ CoDel FQ CoDel has a solution for this problem by providing a fair queuing mechanism.

PFIFO Another alternative is Priority FIFO (PFIFO). As the name suggests it combines the concepts of a basic FIFO queue with priorities. This can reduce delay

5.1.3 Interactive vs PEP

Our first experiment consists of having a interactive flow (100 byte UDP packets at 2kBps) competing with a file transfer, one default end to end and one through our PEP. Highlighting some of the initial differences between an end to end connection and a PEP.

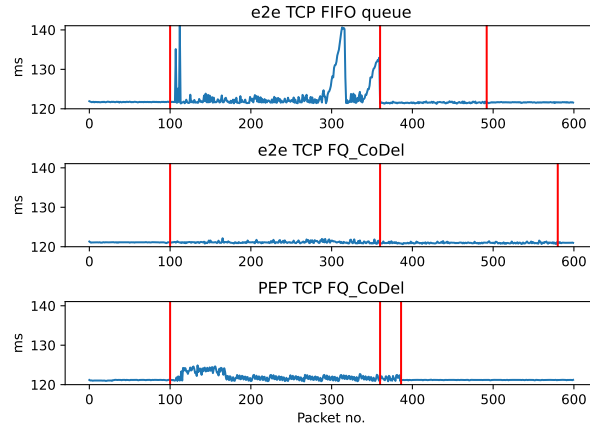


Figure 5.1: Interactive UDP traffic

5.1.4 PEP vs E2E tests

The first test consists of evaluating our PEP against a default TCP end to end (E2E) connection, while also highlighting the difference a packet scheduler can make. Fig. 5.1 shows the results, the red lines represent important events in the timeline. The first line represents the start of a file transfer, the second shows a bandwidth change from 10mbit to 75mbit, while the last line shows when the file transfer finished.

The first timeline shows an end to end TCP transfer using a FIFO queue, in this case BFIFO.

On the second timeline we see the behavior of an end to end connection using FQ CoDel. MORE The last timeline shows the PEP using FQ CoDel. Looking at the graph we can see two important differences between the

default end to end TCP and the PEP. Firstly, the PEP has higher latency fluctuations than the E2E connection, shown by the blue lines on the graph.

5.1.5 10x10 tests

To further evaluate the PEP we conducted an experiment where we have 1-10 flows competing, once using the PEP and once end to end. The goal is to see how competing flows affect the behavior of our PEP

5.1.6 Spike

Chapter 6

Conclusion

6.1 Future Work

Bibliography

- [1] S. K. Agrawal and Kapil Sharma. 5g millimeter wave (mmwave) communications. In *2016 3rd International Conference on Computing for Sustainable Global Development (INDIACom)*, page 3630–3634, Mar 2016.
- [2] Olivier Bonaventure, Mohamed Boucadair, Sri Gundavelli, SungHoon Seo, and Benjamin Hesmans. 0-RTT TCP Convert Protocol. RFC 8803, July 2020.
- [3] David Borman. TCP Options and Maximum Segment Size (MSS). RFC 6691, July 2012.
- [4] Kristjon Ciko, Michael Welzl, and Peyman Teymoori. Pep-dna: A performance enhancing proxy for deploying network architectures. In *2021 IEEE 29th International Conference on Network Protocols (ICNP)*, pages 1–6, 2021.
- [5] Wesley Eddy. Transmission control protocol (tcp). Request for Comments RFC 9293, Internet Engineering Task Force, Aug 2022.
- [6] Toke Høiland-Jørgensen, Paul McKeeney, dave.taht@gmail.com, Jim Gettys, and Eric Dumazet. *The Flow Queue CoDel Packet Scheduler and Active Queue Management Algorithm*. Number RFC 8290. Jan 2018.
- [7] Michele Polese, Marco Mezzavilla, Menglei Zhang, Jing Zhu, Sundeep Rangan, Shivendra Panwar, and Michele Zorzi. milliproxy: A tcp proxy architecture for 5g mmwave cellular systems. In *2017 51st Asilomar Conference on Signals, Systems, and Computers*, pages 951–957, 2017.
- [8] Cristian García Ruiz, Antonio Pascual-Iserte, and Olga Muñoz. Analysis of blocking in mmwave cellular systems: Application to relay positioning. *IEEE Transactions on Communications*, 69(2):1329–1342, Feb 2021.
- [9] W. V. Wathsala, Buddhika Siddhisena, and Ajantha S. Athukorale. Next generation proxy servers. In *2008 10th International Conference*

on Advanced Communication Technology, volume 3, pages 2183–2187, 2008.

- [10] Michael Welzl. Network congestion control: Managing internet traffic. *Network Congestion Control: Managing Internet Traffic*, pages 1–263, 05 2006.
- [11] Michael Welzl, Dimitri Papadimitriou, Bob Briscoe, Michael Scharf, and Michael Welzl. Open Research Issues in Internet Congestion Control. RFC 6077, February 2011.