

# Queueing at the Telco Service Edge: Requirements, Challenges and Opportunities

Ralf Kundel  
Multimedia Communications Lab,  
TU Darmstadt, Germany  
ralf.kundel@tu-darmstadt.de

Joerg Wallerich  
Deutsche Telekom Technik GmbH,  
Darmstadt, Germany  
joerg.wallerich@telekom.de

Wilfried Maas  
Deutsche Telekom Technik GmbH,  
Darmstadt, Germany  
wilfried.maas@telekom.de

Leonhard Nobach  
Deutsche Telekom Technik GmbH,  
Darmstadt, Germany  
leonhard.nobach@telekom.de

Boris Koldehofe  
Multimedia Communications Lab,  
TU Darmstadt, Germany  
boris.koldehofe@tu-darmstadt.de

Ralf Steinmetz  
Multimedia Communications Lab,  
TU Darmstadt, Germany  
ralf.steinmetz@tu-darmstadt.de

## ABSTRACT

Internet Service Providers (ISP) are challenged by high cost awareness and increasing requirements at the same time. Open and programmable hardware, driven by the P4-programming language, enables new possibilities regarding packet header processing on switches. Focusing on data center requirements only, these devices currently do not fulfill all packet queueing requirements of large scale ISPs. In this work we give insights in ISP access networks and specify the resulting queueing demands and required features of future data plane hardware. In addition, we present results of a real world measurement that contributes to a deeper understanding of the queueing requirements in access networks. Last, the potential and importance of active queue management and flexible programmable schedulers are highlighted.

## KEYWORDS

Service Edge, QoS, Massive Queueing, Hierarchical Scheduling, Bufferbloat, Operational Experiences

## 1 INTRODUCTION

The main task of Internet Service Providers (ISP), often referred to as Telco, is the creation of Internet connectivity for their customers but also providing other products as telephony and television. Besides consumers, content, service and application providers, mainly conventional Data Centers (DC) and Content Delivery Networks (CDN), can belong to the group of ISP-customers. Nowadays ISP infrastructures all services, including telephony and television, are realized over Ethernet/IP-based networks and therefore the main task of an ISP forwarding packets to and from its customers. Figure 1 depicts a high-level view on a typical ISP topology.

Typically, the link between the consumer residential gateway (RG) and the ISP is the bottleneck of almost all flows due to technical limitations or contractual restrictions. In order to minimize flow completion times, affected by packet loss,

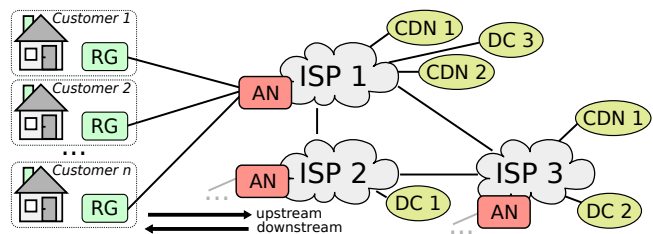


Figure 1: Internet Service Provider topology including customer residential gateways (RG), access networks (AN), content delivery networks (CDN) and traditional data centers (DC).

and enable a higher link utilization, queueing is required in both directions. All packets from the consumer to the ISP, named *upstream*, are queued at the residential gateway. This can be easily done in software as part of the residential gateway as the total upstream bandwidth of one single consumer is in the range of hundreds of Megabyte or even below.

Traffic from the ISP to all connected RGs, named *downstream*, must be queued at the edge of the ISP network, called *service edge*. The focus of this paper will be downstream traffic as it causes the following challenges: 1) Downstream traffic of many customers, which is much more distinct compared to upstream, must be processed and accounted at the same place and time. 2) Complex scheduling restrictions and multiple Quality of Service (QoS) classes must be considered. This paper:

- analyzes the constraints of queueing downstream traffic at the Telco service edge,
- provides real measurement data from access network traffic which assist queue size engineering and
- names requirements for future data plane hardware regarding a better suitability in Telco access networks and higher quality of service.

The outline of this paper is as follows: First, we describe typical access network topologies and from that we derive packet queueing constraints. Second, we present and discuss measurement results from a typical access network. Third, queueing requirements, based on the previous sections, are given and finally related work is discussed.

## 2 TELCO ACCESS NETWORKS

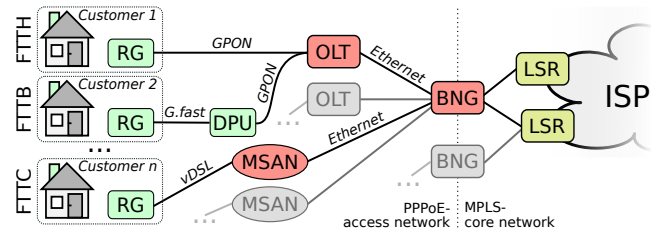
Due to heterogeneous environmental conditions and grown network infrastructures, Internet access is typically realized with several technologies at the same time by a single ISP. Figure 2 depicts the access network topology of *Fiber to the Home (FTTH)*, *Fiber to the Building (FTTB)* and *Fiber to the Curb (FTTC)*. Further access technologies, e. g., DOCSIS or aDSL, will not be considered in detail, but we assume similar challenges and characteristics. In case of the three given access scenarios a tunneling protocol, e. g., PPPoE, is used between each residential gateway (RG) and the broadband network gateway (BNG). The BNG terminates the customer tunnel and is responsible for handling upstream and downstream traffic, including packet header processing, authentication, authorization and traffic shaping. In our previous work [7] we have discussed the PPPoE session termination and have shown the feasibility of this on programmable data plane hardware with the description language P4 [3]. In the following we will describe the differences between these technologies:

**FTTH:** The customer residential gateway is connected via a passive optical network, e. g., GPON or XGSPON, with the optical line terminal (OLT). In contrast to Ethernet connections, a single fiber from the OLT is split optically to  $n$  customers, dependent on the technology up to 64, and therefore a shared resource. The OLT device has a common Ethernet point-to-point connection to the BNG.

**FTTB:** In contrast to FTTH the GPON fiber connection is terminated by a distribution point unit (DPU) and not by the residential gateway. This DPU is typically located in the basement of a building and multiple customer RGs can be connected via G.fast copper wires to it. Advantage of this technology is only little construction works inside houses as existing copper cables can be used for the last meters.

**FTTC:** In case of vDSL Internet access a copper based connection between the residential gateway and the multi service access node (MSAN), formerly known as DSLAM, is used. Compared to FTTB, the distances are higher, typically in the range of hundreds of meters, and by that the physical link bandwidth is more restricted. The MSAN is, like to the OLT, connected to the BNG.

This access network topology can be seen as a tree with the BNG as root node. Downstream traffic, from the ISP to the customer, will be shaped to the agreed bandwidth of the



**Figure 2: Heterogeneous ISP access network topology. Customer PPPoE sessions are terminated at the Broadband Network Gateway (BNG) which is the border to the MPLS routed ISP core network.**

customer at the BNG and all delivered packets are accounted. The following access network, consisting of aggregation devices as a MSAN, OLT or DPU, has many shared resources and by that bandwidth bottlenecks can occur. For example, the sum of all theoretical customer bandwidths at a single MSAN is higher than the shared link to the BNG. Under normal circumstances in a well-designed access network this limit will never be reached, however in such cases a fair and QoS-aware sharing of the bandwidth must be ensured. The packet queueing and hierarchical scheduling, which is aware of the access network limitations, must be performed in the BNG before counting the downstream traffic.

In contrast to most data center networks, ISP access networks are challenged by the previously introduced hierarchical access network, many customers and different QoS-classes which will be discussed in the following sections:

### 2.1 Hierarchical Scheduling

The decision, which packet should be forwarded next by the access network, depends on many parameters:

- **Bandwidth limit:** Each customer has a maximum bandwidth which can not be exceeded.
- **Customer separation:** Interference of different customers must be excluded or at least minimized. E. g., *Customer 1* utilizes his downstream bandwidth fully and *Customer 2* should not suffer under higher latencies by that, implying separated queues.
- **QoS awareness:** If an ISP offers multiple service classes the scheduler should prioritize the flows accordingly. This can be either done within the traffic classes of a single customer only or under consideration of other customers as well.
- **Many queues:** Up to 35.000 customers [9] should be terminated on a single BNG and as each customer requires at least one queue many queues are required.
- **Zero loss:** Packets are not allowed to be dropped somewhere in the access network between BNG and RG

as they are already accounted by the BNG. For that, hierarchical scheduling policies are required.

Considering these criteria, a scheduling decision must be done within hundreds of  $ns$ . Assuming an average packet size of  $1000Byte$  and  $100Gbit/s$  link speed, only  $80ns$  are left between packets. This implies the use of hardware schedulers as software based approaches can not guarantee that.

## 2.2 Quality of Service

In case of a QoS-sensitive ISP, each customer of an ISP utilizes multiple QoS-classes. Table 1 names and describes 6 exemplary traffic classes. *Note:* This arrangement is only an example and will differ from ISP to ISP. E. g., 3 classes can be sufficient whereas some ISPs might use 8 classes.

Each customer  $i$  has a maximum bandwidth  $C_{max,i}$ , shared by all traffic classes. E. g., if the bandwidth is only requested by BE traffic,  $C_{max,i}$  can be used for that. But as soon as another flow belonging to a higher priority class, e. g., IPTV, requests bandwidth by sending packets, the maximum available bandwidth of BE will be reduced.

Due to the hierarchical access network and hierarchical scheduling a priority awareness on shared resources is possible as well. For instance, if many subscribers utilize an OLT in downstream direction fully, phone calls of other subscribers should still be possible without disturbance or additional delays. If and how this should be realized in detail depends strongly on the ISP policies and maybe on national regulations. If a flow of *Customer A*, declared as priority traffic ( $> BE$ ), is able to take precedence over flows of *Customer B* misuse might occur. In order to prevent this, two policies can be applied: (1) matching on source and destination addresses allows the verification of prioritized traffic, e. g., multicast IPTV, and (2) strong rate limiting on different traffic classes, e. g., a few hundreds  $kbyte/s$  for VoIP which is enough for multiple parallel voice calls. Besides misuse, this traffic class aware rate limiting guard against misbehavior of any kind and reason.

These priorities require the use of one FIFO-queue for each QoS-class per customer. Using the same FIFO-queue for all classes would affect in higher latencies, equal for all traffic classes, and head of line blocking due to different bandwidth limits. Push-In-First-Out (PIFO) queues might be an enabler for single queue approaches and will be discussed later in Section 5. For similar reasons, each customer requires an own set of queues as otherwise the customer separation can not be guaranteed.

## 2.3 Buffer Sizing

Queue sizes in general and in this special scenario depend on many influencing factors which causes a major challenge for ISPs due to heterogeneous traffic patterns. Appenzeller

Class	Description
BE <sub>low</sub> (0)	Best Effort low priority (BE <sub>low</sub> ) can be used for "WIFI to go" products where the residential gateway is a public WIFI access point at the same time.
BE (1)	Best Effort (BE) traffic includes almost all traffic of a customer, paying for this Internet connection.
LD/LL (2/3)	Low Delay (LD) and Low Loss (LL) are used for enterprise applications of business customers only. For residential access they are typically not used.
MC (4)	Multicast (MC) is used for television products of ISPs. Duplication is realized at the BNG and by that data transfer in the core network is strongly reduced. <i>Note:</i> This includes only live TV offered by the ISP, not (third-party) on-demand video streaming.
VoIP (5)	Many ISPs offer Voice over IP (VoIP) services to their customers as well which require a very low latency and loss rates. <i>Note:</i> All application layer approaches, e. g., Skype or WhatsApp, are handled as BE traffic.
Ctrl (6)	This traffic class is only used for controlling the ISP network including access network and RGs.

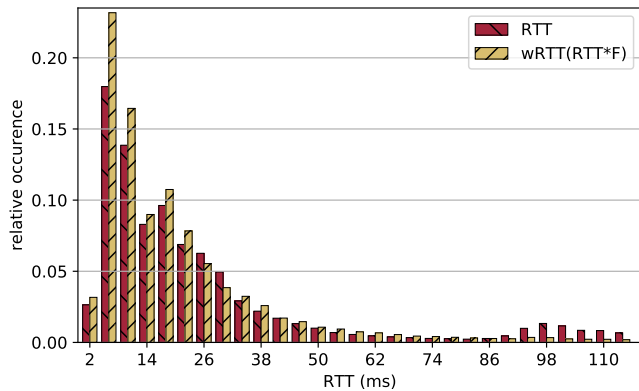
**Table 1: Typical ISP QoS-classes, ascending priorities.**

*et al.* established in 2004 the rule of thumb  $B = \frac{RTT * C}{\sqrt{n}}$  [1] for TCP.  $B$  describes the optimal queue depth in *byte* whereas  $RTT$  describes the (average) round-trip time of the queued flow(s).  $C$  is the queue service rate in *byte/s* and  $n$  describes the number of flows.

In later related work [4][8] similar formulas are presented. Consensus of all of them is, that a higher RTT and link speed requires a bigger queue whereas more parallel congestion controlled flows reduce the optimal queue size. One major influencing factor, not considered by this formula, is the used congestion control mechanism, e. g., CUBIC, Reno or BBR. As this congestion control mechanism is not transparent for network middleboxes, including the BNG, an optimization on that is not possible.

On the one hand, if the chosen queue size is smaller than the optimal one, the maximum bandwidth will not be reached due to a saw-tooth like bandwidth oscillation. On the other hand, most TCP congestion controls tend to fill queues until packet loss occur. By that, a large queue causes large delays even if the optimal queue size is much smaller. This phenomenon is called *Bufferbloat* [5] and mostly unavoidable for fixed sized queues with varying flow RTTs and flow numbers. Queue management algorithms which are supported by current hardware switches are typically taildrop, RED, WRED or similar ones and therefore not able to tackle the Bufferbloat phenomenon.

Active Queue Management (AQM) algorithms address this issue by dynamically controlling queue sizes by either dropping or marking packets. In previous work have shown that AQMs, e. g., CoDel, can be realized on P4-programmable data planes within limitations [6] as well. However, current state-of-the-art switches with a fixed behavior don't support



**Figure 3: RTT distribution for all  $N \approx 70.000.000$  captured TCP flows and weighted by the byte length  $F$  of each flow. Covering 89.44% of RTT and 96.38% of wRTT.**

such algorithms and as P4 is also not made for scheduling, specially complex algorithms considering QoS classes and hierarchies, a production deployment has many hurdles. Indeed, current first approaches of applying advanced active queue management approaches in residential gateways(RG) for upstream traffic exist which can be applied by software updates.

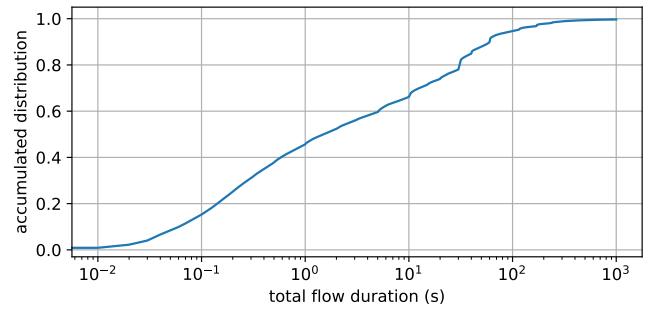
In contrast to best effort traffic, priority classes, such as VoIP or MC, typically have no congestion control and little bandwidth making small queues ( $1ms - 10ms$ ) sufficient.

### 3 MEASUREMENT RESULTS

Designing downstream queues at the BNG makes is necessary to understand the traffic which should be queued. For that, we provide up-to-date measurement results, captured between the first LSR core router and the BNG, which contains all downstream packets and by that the access network traffic. In total, we captured a 24 hours trace, containing around 54 billion packets and 70 million TCP flows. The following measurement results represent the average over a 24 hour trace which has a utilization peak in the evening as discussed in related work before [2]. Note: The shown distributions are almost independent on the daily and weekly periodicity, but we assume a slow transformation over years due to Internet topology changes, e. g., edge clouds and CDNs.

As discussed in the previous section, the optimal buffer size for congestion controlled flows depends on a typically fixed link speed (C), flow RTT and number of flows (n). Thus, we will focus on these metrics in the following.

Table 2 depicts the distribution of all captured unicast packets between TCP, QUIC and UDP (non QUIC). QUIC traffic was determined based on the Port 443 and is not included in the UDP statistics. As you can see, the sum of TCP and QUIC traffic is around 93% percent and therefore congestion controlled. Multicast traffic is not considered as the



**Figure 4: Accumulated TCP flow duration distribution up to 1000s, covering 99.67% of all flows.**

duplication is done at the BNG and therefore these results would be not meaningful. Note that a continuous multicast stream requires only tiny buffers as this is burst-free.

#### 3.1 RTT Distribution

As TCP is the main part in the traffic mix and can be investigated very well, we will focus next on the RTT distribution of all TCP connections as depicted Figure 3. Each TCP flow RTT was determined at the initial SYN, SYN-ACK, ACK handshake. By this early RTT measurement a falsification by this flow itself, e. g., due to bufferbloat, can be precluded. We observe that the RTT of most flows is below  $50ms$ . A very likely assumption is that a lot of connections are established to CDN servers, located close to the customers. In order to consider the RTT in relation to the transmitted data, we introduce the weighted RTT of flow  $i$ ,  $wRTT_i = RTT_i * F_i$  with  $F_i$  representing to flow length in bytes. The comparison of the RTT and wRTT distribution illustrates that the total delays of these “power flows” are even shorter in average. These power flows benefit from optimal buffer sizes, as they are bandwidth hungry and tend to cause bufferbloat at the same time. Note: This RTT distribution is based on TCP flows only. However, for QUIC we assume very similar RTTs.

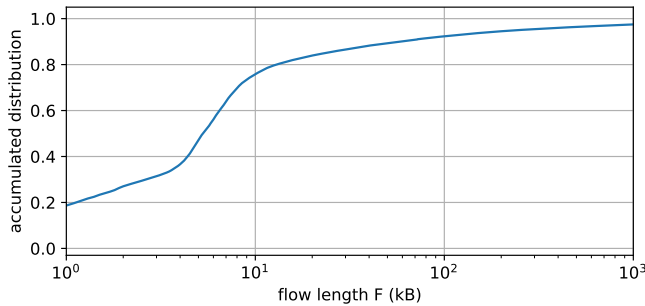
#### 3.2 Flow Characteristics

Figure 4 shows the accumulated flow duration distribution on a logarithmic time scale for all observed TCP connections. Around 50% of the flows has a duration of 1s or less and therefore will not be able to bloat a buffer.

In Figure 5 the flow distribution based on the flow length (F) is given. Only 20% of the TCP flows have a total length of 10KB or more. We assume that this 80% of the flows are tiny

	TCP	QUIC	UDP	other
share:	84.4%	8.8%	6.4%	0.4%
avg. pkt. size [byte]:	1347	1319	832	721

**Table 2: Unicast traffic distribution.**



**Figure 5: Accumulated flow length (F) distribution over all TCP flows, covering 97.48% of all flows.**

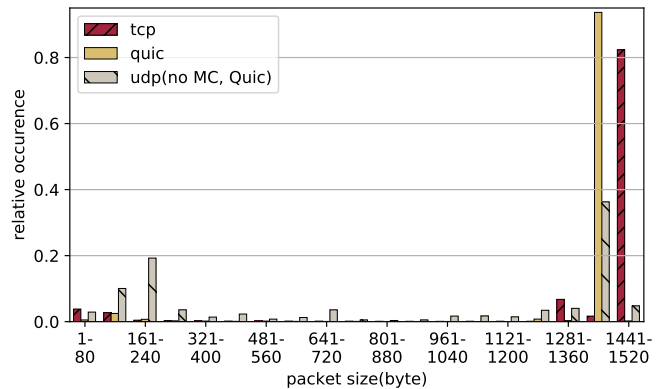
requests, e. g., retrieving many css or javascript files as part of loading a website.

Combined with the previous duration measurement results we can say that there are many flows with a short duration and little transmitted data. The flow completion time of these tiny flows is mainly depended on the RTT and not on the available bandwidth. A higher RTT would cause a slower congestion window increase and by that a smaller bandwidth utilization. Therefore, they are vulnerable to bufferbloat but at the same time will not cause this phenomenon. Note: Although these flows have a huge total number, their share on the total bandwidth is quite low as only a few big flows causes most bandwidth consumption.

Another interesting observation is the distribution of packet sizes, as depicted in Figure 6 and Table 2. Note: The given values does not include packet header overhead caused by the ISP, e. g., a MPLS header stack of the core network. Whereas TCP utilizes the Ethernet MTU fully, QUIC implementations stay slightly below. Furthermore, it is visible that most packets are very close to the Ethernet MTU of 1500 byte and only a few of them are smaller. A second cluster can be detected in the range of  $\leq 240$  byte. This information is quite helpful for designing switches and other network components, buffering packets. A memory allocation of MTU bytes for each packet will cause only little unused memory overhead compared to a per packet memory allocation considering the size of each packet. In addition, it might be useful to allocate small memory ranges as well for packets of  $\leq 240$  byte if memory is a narrow resource.

### 3.3 Quality of Service

Last, in Table 3 the distribution over different QoS classes is depicted. It is obvious that almost all traffic belongs to the class of best effort (BE). This can be ensured by strong bandwidth limits for all priority classes. By that, designing a hierarchical scheduler will be slightly simplified as an over-utilization of shared resources can never occur by non best effort traffic only and by that a guaranteed forwarding can



**Figure 6: Distribution of packet length for unicast traffic, categorized by TCP, QUIC and other UDP traffic. Packets are grouped in blocks of 80 bytes, 1-80 byte, 81-160 byte, ... .**

Class:	BE	LD/LL	VoIP	Ctrl
share:	99.82%	0.03%	0.14%	0.01%
avg. pkt. size [byte]:	1314	382	200	886

**Table 3: Unicast traffic distribution on QoS-classes.**

be ensured for all priority classes. However, this distribution represents only a current snapshot and will might change strongly in future due to more QoS-sensitive applications, e. g., IoT.

## 4 SCHEDULER REQUIREMENTS

Based on the presented requirements regarding queueing and scheduling at the service edge of ISPs and the measurement results, we summarize the following requirements. As the number of customers and QoS-classes is non-uniform for different ISPs, the number of required queues, which is the product of these two values, might strongly differ. Therefore, only a realistic upper bound is given. A priority aware, hierarchical scheduler will be required as well in order to guarantee no over-utilization of the access network. However, the number of hierarchical layers will differ from ISP to ISP as well as the concrete scheduling algorithm. Based on current network topologies and related work [7] we assume a maximum total downstream bandwidth per BNG of 100 – 500Gbit/s to be realistic.

If AQM algorithms [6] can be applied on future service edges in hardware, the required queue memory can be reduced as in most cases the RTTs are quite low. Combining the rule of thumb  $B = \frac{RTT * C}{\sqrt{n}}$  [1] with an average RTT of 20ms,  $C = 100Gbit/s$  downstream bandwidth and  $n = 2000$  parallel flows per BNG, the total required memory space is 56MB. This is the borderline of nowadays on-chip packet buffers

<b>Number of Customers:</b>	5.000 – 35.000
<b>Number of QoS-Classes:</b>	3 – 8
<b>Number of Queues:</b>	≤ 280.000
<b>Scheduling:</b>	Priority/QoS aware, 3-6 hierarchical layers
<b>Queue Sizes:</b>	AQM, 1ms – 100ms (QoS-class dependent)
<b>Downstream Bandwidth:</b>	100Gbit/s – 500Gbit/s

**Table 4: BNG queueing and scheduling requirements.**

on network switches and consequently no external memory is required. Considering that customers which do not utilize their bandwidth fully require almost no queue memory as packets can be forwarded immediately. Last, we want to note the possible support of Explicit Congestion Notification (ECN) at the service edge. With future programmable data plane hardware, including powerful programmable schedulers, advanced queueing approaches like *ECN-LAS* can be deployed with ease and provide a benefit to customers.

## 5 RELATED WORK

Sivaraman *et al.* have introduced an alternative, called PIFO, to classical FIFO-queues which enables the insertion of a packet at any position in the queue, e. g., based on its QoS class [10]. One single *Push-In-First-Out (PIFO)* can be used per customer if strict priorities between the QoS classes should be applied. Furthermore, the authors introduced a description language for flexible scheduler programming which can be compiled to a generic hardware.

Other approaches of combining many flows in a single queue exists as well. E. g., the work *Urgency Based Scheduler* [11] has shown that many flows with equal bandwidth limits can be combined in a single queue with bandwidth guarantees for each flow. However, customer separation in terms of latency interference can not be guaranteed by that.

## 6 CONCLUSION

Queueing at the service edge of ISPs is a very important but also challenging function as by that bandwidth utilization, fairness and QoS-awareness must be ensured. In this work we have analyzed the requirements of queueing downstream traffic regarding queue sizes and scheduling algorithms, challenged by hierarchical scheduling constraints, many queues, different QoS-classes and heterogeneous traffic patterns. For a better understanding of the variables influencing the optimal queue size, we provide real-world measurement data of downstream traffic. These results show that end-to-end latency can be reduced in access networks by methods to determine the optimal queue size, e. g., by AQM algorithms. Last, we analyzed and determined the parameters of a fully-fledged, carrier-grade queueing implementation. In future

work we will investigate the feasibility of these requirements on programmable standard hardware, mainly FPGAs.

## ACKNOWLEDGMENT

This work has been supported by Deutsche Telekom through the Dynamic Networks 8 project, and in parts by the German Research Foundation (DFG) as part of the project C2 within the Collaborative Research Center (CRC) 1053 MAKI. Furthermore, we thank our colleagues for their valuable input and feedback.

## REFERENCES

- [1] Guido Appenzeller, Isaac Keslassy, and Nick McKeown. 2004. Sizing Router Buffers. In *Proceedings of the 2004 Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications (SIGCOMM '04)*. ACM, New York, NY, USA, 281–292.
- [2] Andreas Betker, Inken Gamrath, Dirk Kosiankowski, Christoph Lange, Heiko Lehmann, Frank Pfeuffer, Felix Simon, and Axel Werner. 2014. Comprehensive Topology and Traffic Model of a Nationwide Telecommunication Network. *J. Opt. Commun. Netw.* 6, 11 (Nov 2014), 1038–1047.
- [3] Pat Bosshart, Dan Daly, Glen Gibb, Martin Izzard, Nick McKeown, Jennifer Rexford, Cole Schlesinger, Dan Talayco, Amin Vahdat, George Varghese, and David Walker. 2014. P4: Programming Protocol-independent Packet Processors. *SIGCOMM Comput. Commun. Rev.* 44, 3 (July 2014), 87–95.
- [4] Amogh Dhamdhere and Constantine Dovrolis. 2006. Open Issues in Router Buffer Sizing. *SIGCOMM Comput. Commun. Rev.* 36, 1 (Jan. 2006), 87–92.
- [5] Jim Gettys and Kathleen Nichols. 2012. Bufferbloat: Dark Buffers in the Internet. *Commun. ACM* 55, 1 (Jan. 2012), 57–65.
- [6] R. Kundel, J. Blendin, T. Viernickel, B. Koldehofe, and R. Steinmetz. 2018. P4-CoDel: Active Queue Management in Programmable Data Planes. In *2018 IEEE Conference on Network Function Virtualization and Software Defined Networks (NFV-SDN)*. 1–4.
- [7] Ralf Kundel, Leonhard Nobach, Jeremias Blendin, Hans-Joerg Kolbe, Georg Schyguda, Vladimir Gurevich, Boris Koldehofe, and Ralf Steinmetz. 2019. P4-BNG: Central Office Network Functions on Programmable Packet Pipelines. In *15th International Conference on Network and Service Management (NOMS)*. IFIP.
- [8] Wolfram Lautenschlaeger. 2014. A deterministic tcp bandwidth sharing model. *arXiv preprint arXiv:1404.4173* (2014).
- [9] Dan Rodriguez. 2018. Next Generation Central Offices Transform Network Edge with Datacenter Economics, Cloud Flexibility.
- [10] Anirudh Sivaraman, Suvinay Subramanian, Mohammad Alizadeh, Sharad Chole, Shang-Tse Chuang, Anurag Agrawal, Hari Balakrishnan, Tom Edsall, Sachin Katti, and Nick McKeown. 2016. Programmable Packet Scheduling at Line Rate. In *Proceedings of the 2016 ACM SIGCOMM Conference (SIGCOMM '16)*. ACM, New York, NY, USA, 44–57.
- [11] J. Specht and S. Samii. 2016. Urgency-Based Scheduler for Time-Sensitive Switched Ethernet Networks. In *2016 28th Euromicro Conference on Real-Time Systems (ECRTS)*. 75–85.