# A Survey on Recent Advances in Transport Layer Protocols

Michele Polese, *Student Member, IEEE*, Federico Chiariotti, *Member, IEEE*, Elia Bonetto, Filippo Rigotto, Andrea Zanella, *Senior Member, IEEE*, Michele Zorzi, *Fellow, IEEE* 

Abstract—Over the years, the Internet has been enriched with new available communication technologies, for both fixed and mobile networks and devices, exhibiting an impressive growth in terms of performance, with steadily increasing available data rates. The Internet research community has kept trying to evolve the transport layer protocols to match the capabilities of modern networks, in order to fully reap the benefits of the new communication technologies. This paper surveys the main novelties related to transport protocols that have been recently proposed, identifying three main research trends: (i) the evolution of congestion control algorithms, to target optimal performance in challenging scenarios, possibly with the application of machine learning techniques; (ii) the proposal of brand new transport protocols, alternative to the Transmission Control Protocol (TCP) and implemented in the user-space; and (iii) the introduction of multipath capabilities at the transport layer.

#### I. INTRODUCTION

The communication technologies that provide access and backhaul connectivity to the Internet have dramatically changed since the 1980s, when the protocols that are part of today's TCP/IP stack were first introduced [1]-[3]. Traffic generated on mobile devices is expected to exceed desktop and server traffic by 2021 [4]. New communication standards are being proposed and launched to market every few years. Driven by the increase in web and multimedia traffic demand, mobile and fixed networks are rapidly evolving. 3GPP NR [5], [6] will bring ultra-high data rates with low latency to future 5G devices, and, similarly, the IEEE 802.11 standard will target ultra-dense deployments [7] and multi-gigabit-per-second throughput [8]. Modern devices are capable of connecting to heterogeneous networks [9], and fixed networks are using new optical technologies and Software-Defined Networking (SDN) for unprecedented rates and low latency [10].

The increasing capabilities of the network make new kinds of applications possible; the exponential growth of multimedia or real-time traffic [4] would have been impossible without the recent technological advances. As networks progress towards 5G, new kinds of applications, such as Augmented Reality (AR)<sup>1</sup> and Virtual Reality (VR) or autonomous cooperative

This work was partially supported by the program Supporting Talent in Research@University of Padua: STARS Grants, through the project "Cognition-Based Networks: Building the Next Generation of Wireless Communications Systems Using Learning and Distributed Intelligence," and by the US Army Research Office under Grant no. W911NF1910232: "Towards Intelligent Tactical Ad hoc Networks (TITAN)."

<sup>1</sup>A comprehensive list of acronyms is provided at the end of the paper.

driving, are going to require more from the network and impose ever more stringent Quality of Service (QoS) constraints. This, along with the increasing heterogeneity of the network, makes the role of transport protocols more important, and, at the same time, more challenging. Indeed, the end-to-end performance and the Quality of Experience (QoE) of the users largely depend on the interaction among the applications, the transport layer and the underlying network [11]. In particular, the transport layer, which is responsible for the management of the end-to-end connection over any number of network hops, has to adapt and evolve in order to let users fully benefit from the aforementioned innovations. However, a number of factors prevent new solutions at the transport layer from being widely adopted, and, in recent years, the research community has been forced to cope with these limitations and identify innovative solutions in order to have significant effects on Internet performance.

1

In particular, the deployment of alternative transport protocols, such as the Stream Control Transmission Protocol (SCTP) [12], is slowed down by the widespread use of middleboxes [13], [14], which often drop packets from protocols which are different from the Transmission Control Protocol (TCP) and/or the User Datagram Protocol (UDP) [15]. Moreover, the socket Application Programming Interface (API) (offered by the Operating System kernel and supported by TCP/UDP) is almost universally used [16], thus limiting the interfacing options between application and transport layers to what the API supports. Finally, the most widespread Operating Systems implement the transport functionalities (i.e., TCP and UDP) in the kernel, making the deployment of new solutions difficult. These elements define what is called transport layer ossification [16], a phenomenon which has pushed developers and researchers to only use legacy TCP (for reliable traffic and congestion control), even though it may not be the best performing protocol for the desired use case. TCP has indeed some performance issues in specific scenarios, e.g., on wireless links with high variability [17], [18], Head of Line (HoL) blocking with web traffic [19], and bufferbloat [20].

This survey focuses on three directions in transport layer research (i.e., new transport protocols, congestion control innovations and multipath approaches) that have emerged in the last fifteen years to solve the aforementioned problems. New congestion protocols have been proposed to target low latency and full bandwidth utilization [21], [22]. Novel transport protocols have been discussed by the Internet Engineering Task Force (IETF), with technical novelties such as multipath capabilities, to exploit the multiple interfaces available in

The authors are with the Department of Information Engineering (DEI), University of Padova, Padova, 35131, Italy. Email:{polesemi, chiariot, bonettoe, rigottof, zanella, zorzi}@dei.unipd.it

modern smartphones, computers and servers, and user space implementations, to overcome the ossification that prevents a widespread adoption of novel algorithms at the transport layer. Therefore, in this survey, we first review the main new transport protocols that have been proposed or standardized by the IETF since 2006: we provide a brief overview on SCTP and Datagram Congestion Control Protocol (DCCP) [23], and then delve into a more recent contribution, i.e., the Quick UDP Internet Connections (QUIC) [24] protocol. Then, we review the research on congestion control. We describe both new mechanisms using classic approaches, e.g., Bottleneck Bandwidth and Round-trip propagation time (BBR) and Low Latency (LoLa) for TCP, and some novel proposals for using machine learning techniques for congestion control. Finally, the third trend we discuss in this survey is related to the adoption of multipath solutions at the transport layer, mainly with Multipath TCP (MPTCP) [25], [26], but also with multipath extensions for SCTP [27], QUIC [28], [29], and with the Latency-controlled End-to-End Aggregation Protocol (LEAP) [30].

Our goal in this survey is to offer the interested reader a comprehensive point of view on the up-to-date research on transport protocols, which is lacking in other recent surveys that separately focus on ossification [16], multipath transmissions [31], or congestion control schemes for MPTCP [32] or in data centers [33].

The remainder of the paper is organized as follows. In Sec. II we describe the main issues and limitations of transport protocols in modern networks. Then, in Sec. III, we describe three recently proposed transport protocols, focusing in particular on QUIC. In Sec. IV we report the latest proposals in terms of congestion control algorithms for TCP and other transport protocols, classifying them as loss-, delay- or capacity-based. We also discuss hybrid mechanisms, machine-learning-based algorithms and cross-layer approaches. In Sec. V we investigate the different approaches to multipath transport protocols recently proposed. Finally, Sec. VI concludes the paper and summarizes the main future research directions. A comprehensive list of acronyms is also provided at the end of the paper.

## II. TRANSPORT LAYER LIMITATIONS IN MODERN NETWORKS

As mentioned in Sec. I, transport layer protocols have an end-to-end view of the connection: they do not consider each individual hop, but only a single logical link between the two endpoints. For this reason, the most important link in the connection is the slowest one, that is the so-called bottleneck. The service provided by TCP, the *de facto* standard transport protocol of the modern Internet, can then be roughly modeled as a single pipe with the capacity of the bottleneck link and a certain Round Trip Time (RTT), i.e., the time from the instant the sender transmits a packet to the instant it receives the corresponding acknowledgment.

However, the particular features of the individual links and the behavior of lower layers do influence TCP's behavior, as well as that of other transport protocols: several properties of the links composing the end-to-end connection (e.g.,

latency, packet loss, buffer state and size, and volatility of the capacity) can affect the transport layer performance [34]. In order to overcome this problem, more and more complex congestion control mechanisms have been proposed, going far beyond the original simple Additive Increase Multiplicative Decrease (AIMD) principle. We will describe the main design philosophies in greater detail in Sec. IV. In fact, researchers are questioning TCP's extensive use as a one-size-fits-all solution, because all these factors can cause performance issues that are becoming more and more apparent. Some of these issues are fundamental problems of the transport layer abstraction, while others depend on the specific features of the protocol and can be avoided by designing the protocol correctly. This section provides a short review of some of the most important issues, while the rest of the paper is dedicated to the discussion of the several solutions that have been proposed to address these challenges.

#### A. Bufferbloat

As we will discuss in depth in Sec. IV, congestion control mechanisms exploit an abstract view of the underlying network to tune the amount of data to be sent. As shown in [20], however, this abstraction in some instances fails to provide accurate information on the links connecting the two hosts and leads to degraded performance. In particular, when large buffers are deployed before a bottleneck in order to prevent packet losses, then loss-based TCP probing mechanisms increase the queue occupancy, thus causing a spike in latency. Moreover, since the currently implemented versions of TCP will keep increasing the sending rate until the first packet loss, they will often overshoot the capacity of the channel, increasing congestion and causing multiple retransmissions when the queue is eventually filled. Other protocols such as QUIC, SCTP, and DCCP face the same issue, since it is a fundamental problem of congestion control with large buffers and not a protocol-specific problem.

This phenomenon, known as *bufferbloat*, degrades the QoS of applications, in particular when video or file transfer flows share the buffer with web browsing flows, and it has worsened in recent years mainly due to loss-preventing design strategies that place large buffers in front of low capacity access links (either wired or wireless) [35].

The research in this area aims at solving this issue with local Active Queue Management (AQM) techniques or end-to-end flow and congestion control for transport protocols.

The problem, while not being hard to detect, is hard to solve without a significant overhead cost [36]. The congestion control protocols at the endpoints might also be different, making AQM more complex; this information is often not even available to routers, which might not be able to predict the consequences of discarding a packet on congestion, making the algorithms extremely complex and sensitive to the parameter settings.

AQM is not a recent idea: numerous techniques were proposed, such as Random Early Dropping (RED), and we refer the reader to the extensive literature cited in [36], [37] for further information on this subject. Despite these many proposals, they have encountered limited adoption, partly due to the above-mentioned issue of parameter tuning and the computational cost of the algorithms [35]. More recently, new and easier ways to tune and deploy AQM schemes, such as Controlled Delay Management (CoDel) [38], have been adopted in several commercial products. CoDel is an AQM algorithm that limits the buffering latency by monitoring the queueing delay D in an interval (typically of 100 ms) and dropping packets when the minimum value of D is larger than 5 ms. Nonetheless, as we will discuss in Sec. II-D, the bufferbloat issue remains relevant in the wireless domain (e.g., at mmWave frequencies [18]).

## B. The Incast issue

Data centers are restricted areas containing servers and systems monitoring server's activity, web traffic and performance. The data exchange between servers generally relies on APIs based on the HyperText Transfer Protocol (HTTP), making TCP a widely used transport protocol in data centers. Some activities, such as virtual machine migrations, also generate a high volume of traffic between servers. Therefore, the links in a data center generally have high bandwidth and low latency and delay, while switches have small buffers [39], contrary to what usually happens in access links, as mentioned in the previous section.

Cloud computing frameworks are also widely deployed in large data centers and generate very high traffic loads. For example, MapReduce (which uses a partition/aggregation design pattern) [40] or PageRank (used for web search) [41] often involve many-to-one traffic patterns, where multiple workers send data simultaneously to a single aggregator node, as shown in Fig. 1. In this many-to-one scenario, if all the multiple incast flows go through a single switch, its buffer might be insufficient, leading to congestion. The TCP loss recovery mechanism will then become less efficient, triggering multiple timeouts and causing throughput collapse and long delays [42].

Many attempts have been made to analyze and solve this problem, called the Incast issue, that degrades network performance and user experience [43]. Detailed throughput estimation analysis can be found in [42] and [44]. Solutions may be classified into four categories, as mentioned in [42]: (i) system parameters adjustments, like disabling slow start to avoid massive and sudden buffer overflows that cause retransmission timeout; (ii) enhanced in-network and clientside algorithm design, to reduce waste of bandwidth, minimize retransmission timeouts and the number of packet losses, and to improve the quick recovery of lost packets [39]; (iii) replacement of loss-based congestion control algorithms with better implementations that adjust the congestion window size according to the delay measured from RTTs, like Vegas; and (iv) design of completely new algorithms for this particular environment, like Data Center TCP (DCTCP) [45], that uses Explicit Congestion Notification (ECN) to provide windowbased control methods, or IATCP [46], a rate-based approach that counts the total number of packets injected to constantly meet the Bandwidth-Delay Product (BDP) of the network [42].



Figure 1: A typical scenario in which the TCP Incast problem arises.

Additional information regarding the transport protocol control in data centers can be found in [47].

#### C. Latency and Head of Line Blocking

HoL is a phenomenon that happens when two or more independent data flows share the same TCP connection, as for example with web traffic over TCP: since this transport protocol interprets the data it receives as a single continuous stream of bytes, and it requires in-order delivery, one missing packet delays all subsequent packets for any other flows, causing significant delays. HoL is a well-known problem of any protocol that requires in-order delivery, and the obvious solution of opening one connection for every data flow suffers from significant overhead in connection setup and error recovery. Moreover, with this option, the congestion control is less stable, since each connection performs it independently [19].

Web traffic is a textbook example of the HoL problem: usually, web pages contain several objects, such as text, images, media and third-party scripts; when a client requests a page to the server, each of these objects is downloaded with a single HTTP GET request, but they do not need to be displayed at the same time. HTTP/1.1 did not allow multiplexing, so the client was forced to open one TCP connection for every object, with the issues described above. Version 2 of the protocol, introduced in 2015 as RFC 7540 [56], was supposed to solve this issue by using a single TCP connection to handle all the requests, with significant page load time reductions. However, this advantage is nullified by HoL in lossy networks: given that all the packets for multiple HTTP 2.0 streams are multiplexed over the same TCP connection, which requires inorder delivery, then the loss of a single packet would halt the reception on all streams, and not just on the one interested by the loss, as reported in Fig. 2. For example, measurements in cellular networks show that HTTP/2 does not have a significant performance advantage over HTTP/1.1 [19].

In principle, in-order delivery is not necessary for reliability, but transport protocols such as TCP require both, and although a proposal exists to extend TCP to allow out-of-order transmission for multimedia services to avoid the HoL problem [57], it has not been widely adopted. Fig. 2 shows that protocols that support multiple data streams such as SCTP or QUIC, which will be described in Sec. III, can overcome the HoL problem because there is no single line to block: each stream has its



Figure 2: The HoL problem with TCP and HTTP/2, solved by QUIC thanks to the support of multiple independent streams per connection.

own buffer and in-order delivery is ensured on a per stream basis [58].

The multipath scenario is another case in which TCP suffers heavily from the HoL problem; we will describe the issue in detail in Sec. V.

#### D. Performance on Wireless Channels

The performance of TCP on wireless channels has been investigated since the introduction of the first commercial wireless services in the 1990s [59]. Traditional congestion control schemes generally react to packet loss by reducing the sending rate, assuming that it is caused by congestion. In wireless links, however, packets may be lost because of degraded channel quality. Therefore, the abstraction of the overall connection is poor, and the end-to-end performance suffers. Common solutions include the provisioning of retransmission and protection mechanisms on the wireless link [60] (even though this increases the delay variability), and innetwork performance-enhancing proxies (e.g., to split the connection and use a different congestion control algorithm on the wireless link) [61]. Several techniques have been proposed throughout the years for different technologies, e.g., LTE [62] and Wi-Fi [63]. We refer the reader to [62], [64] for a more extensive discussion on the performance of TCP on wireless links.

4

Recently, the adoption of the mmWave frequencies for the next generation of cellular networks [65] has sparked a renewed interest in the performance of TCP on wireless links. In particular, in this wireless medium the variability of the channel is much higher than at sub-6 GHz frequencies, given the sensitivity to blockage from common materials, and the directionality of the communications. These limitations not only impact the design of the lower layers of the protocol stack, but also affect the transport layer performance, as discussed in [66], [67]. TCP's control loop is indeed too slow to properly react to the dynamics of the channel quality and offered data rate that affect mmWave links. For example, even though bufferbloat is also present in sub-6 GHz cellular networks [68], this phenomenon is much more disruptive at mmWaves [18], where large buffers are needed to protect from temporary but sudden variations of the link capacity due to, e.g., Line of Sight (LOS) to Non Line of Sight (NLOS) transitions [66]. Moreover, it has been shown that TCP makes a suboptimal use of the very high mmWave data rates, since it takes a long time to reach full capacity after a loss or at the beginning of a connection [18], [66]. Since the issue is inherent in the congestion control logic, any protocol implementing congestion control will need to work efficiently also in a wireless scenario.

#### **III. RECENT TRANSPORT PROTOCOLS**

In this section, we will focus on the description of three recently proposed transport protocols: QUIC, SCTP and DCCP. In particular, we will describe the main features of each protocol, and the novelties introduced with respect to TCP. Table I summarizes the main characteristics of these protocols, together with the relevant IETF documents and some details of their implementations. The discussion of the congestion control mechanisms that can be used by these protocols and the recent multipath extensions will be presented in later subsections.

Table I: Summary of the main features of the three recent transport protocols reviewed in this paper (SCTP, QUIC and DCCP). We also include a review of the legacy TCP and UDP protocols as a comparison.

Protocol	RFC or Drafts	Features	Implementations	Multipath extensions
QUIC <sup>†</sup>	[24], [48]	Byte stream, flow and congestion con- trol, reliability (configurable), multi- streaming, integration with TLS, zero- RTT connection establishment	User space, multiple libraries are avail- able [49]–[52]	Multipath QUIC, described in Sec. V-B
SCTP	[12]	Byte stream, flow and congestion con- trol, reliability (configurable), multi- streaming, multi-homing	Operating system kernel	CMT-SCTP described in Sec. V-B
DCCP	[23], [53], [54]	Datagram based, congestion control	Operating system kernel	MP-DCCP, described in Sec. V-B
TCP	[1], [55]	Byte stream, flow and congestion con- trol, reliability	Operating system kernel	MPTCP, described in Sec. V-A
UDP	[3]	Datagram based, best effort	Operating system kernel	No

<sup>†</sup>QUIC is implemented on top of UDP, but provides transport layer functionalities.





Figure 3: Protocol stack with QUIC and TCP, as described in [58].

#### A. Quick UDP Internet Connections (QUIC)

QUIC is an application layer transport mechanism [69] introduced by Google in 2013 and currently considered for standardization by the IETF [24], [48].<sup>2</sup> QUIC is not implemented in the kernel, but distributed either as a user-space library or in the application itself, and relies on UDP as the underlying transport protocol. The main motivations related to the introduction of QUIC are discussed in [70]. While several alternatives to TCP have been recently proposed at the transport layer, their adoption is not widespread, due to the ossification caused by the presence of middleboxes in the network that drop packets from unknown protocols, i.e., protocols different from TCP and/or UDP.<sup>3</sup> OUIC encapsulates its packets into UDP datagrams, thus avoiding most of the incompatibilities with middleboxes. Secondly, as we mentioned in Sec. I, TCP is implemented in the kernel, and any change or new congestion control algorithm requires an operating system update, which is not always feasible. Finally, in terms of protocol design, QUIC aims at solving the HoL issue of TCP, described in Sec. II-C, and at reducing the handshake delay at the beginning of a connection.

QUIC incorporates (i) some of the TCP features, including the acknowledgment mechanism, congestion control and loss recovery; (ii) the key negotiation of TLS 1.3, requiring an allencrypted connection; and (iii) features of HTTP/2, like multistreaming [58]. Fig. 3 shows the difference in the protocol stack for a combination of HTTP/2 and TCP and HTTP/2 and QUIC. The combination of these elements allows QUIC to optimize several key areas. For example, the integration of transport- and cryptographic-related handshakes in the same messages makes it possible to reduce the time for the initial handshake. Furthermore, packet header encryption makes packet inspection and tracking harder for middleboxes. Moreover, it improves TCP's loss recovery by explicitly distinguishing in the ACKs between lost and out-of-order packets. A single connection is composed of multiple independent



Figure 4: Connection establishment procedures; the blue arrows represent the TCP or QUIC handshake, the black ones the cryptographic handshake. The red arrow is the first data transfer.

streams (which can be mapped to HTTP/2 streams) so that the loss of a packet of a single stream does not block the others, effectively mitigating the HoL issue. Finally, additional identifiers are introduced to label a connection, which can then be maintained even in case of changes in the Internet Protocol (IP) addresses of the endpoints.

1) Connection Establishment: The connection establishment is one of the main novelties introduced in QUIC, especially with the 0-RTT option. TCP requires at least one RTT for its handshake when using TLS 1.3 [71], and one or two additional RTTs for the cryptographic handshake when using the older 1.2 version [58].

QUIC enforces the use of at least TLS 1.3 and, in case of a first-time establishment, it carries over all the relevant QUIC transport setup parameters in the first packet, a *Client Hello* message, so that a single RTT is sufficient. The parameters negotiated during the first connection and the server's Diffe-Hellman value, used to calculate the encryption key, are stored at the client. In case of subsequent connections, this information is sent by the client to the server together with the first data packet, and the server uses it to authenticate the client and decrypt the data, thus realizing the 0-RTT handshake shown in Fig. 4.

2) **QUIC Packets:** The use of TLS 1.3 also allows QUIC packets to traverse middleboxes without having their inside information tampered with. In fact, the only bytes that are not encrypted are the UDP header and a portion of the QUIC header, as shown in Fig. 5: this hides key information on QUIC flows from network equipment and operators. After the connection establishment, the only values that are sent in cleartext in the QUIC header are the source and desti-

<sup>&</sup>lt;sup>2</sup>The set of IETF drafts related to QUIC can be found at https://datatracker. ietf.org/wg/quic/documents/.

<sup>&</sup>lt;sup>3</sup>Notice that some firewalls and middleboxes also drop UDP traffic [69].



Figure 5: Example of structure of a QUIC packet with *long header* for IETF QUIC (version 14), adapted from [24]. The red fields are encrypted, while the blue ones (i.e., the packet number field in the header) are protected against casual observation.

nation connection IDs (to allow routing, identification and IP address changes), and the packet type, version number and packet length. The packet number is not protected with the same encryption that is applied to the payload, but has a confidentiality protection against casual observation (e.g., middleboxes on the path). Packet numbers, as in TCP, are different for sending and receiving flows, with the constraint that they must not be reused within the same connection, not even for retransmissions [24].

The packet payload consists in one or more frames. Each frame has a type, like ack, ping or stream, and a series of fields that are type-dependent, such as *stream ID* and *data length* [24]. Therefore, a packet can multiplex multiple streams.

3) OUIC Performance: Given QUIC's promising design goals, its performance has been recently under the spotlight [69], [70], [72], [73], even though the protocol has not yet been completely specified by the IETF. As of today, the QUIC protocol is implemented in Google servers, and, clientside, in Google Chrome and in some mobile applications like YouTube or Google Search. According to [70], in 2017 it accounted for 7% of all Internet traffic, and up to 30% of Google's egress traffic. A more detailed analysis of QUIC deployment is provided in [74], which shows that the number of QUIC-capable endpoints is increasing, mainly thanks to Google and Akamai deployments, while the traffic is almost exclusively covered by Google and its services, reaching a 6.7% share with respect to TCP/HTTPS in September 2017 along the MAWI backbone, as shown in Fig. 6, and a 9.1%share in a mobile ISP. The authors also point out a possible problem in long-term stability due to the frequent and possibly not backward-compatible updates of the protocol.

Another work [70] discusses some metrics measured from Google's live traffic. Fig. 7 reports the reduction in search latency<sup>4</sup> experienced by clients using QUIC with respect to clients using TCP/TLS. The same paper claims that a reduction



6

Figure 6: Percentage of traffic over QUIC vs other protocols on the MAWI backbone during 2017. Reprinted, with permission, from [74].



Figure 7: Percentage of reduction of *Search Latency*. Reprinted, with permission, from [70].

in the re-buffering rates in the order of 18% has been observed for YouTube traffic. Finally, the authors highlight that, with current implementation at the time of writing, the CPU load is on average 3.5 times higher with QUIC than with TCP/TLS [70].

More recently, another work [69] has tested the performance of 13 different versions of QUIC, in desktop and mobile scenarios. With respect to previous studies [72], [73], this paper provides a more thorough investigation, using a more advanced testbed and more recent versions of the protocol. The main findings (using CUBIC as congestion control) confirm that QUIC performs better than TCP/TLS in a desktop environment, but suffers from packet re-ordering issues (which may arise in a mobile environment). With respect to variations in the underlying data rate, QUIC outperforms TCP thanks to the more advanced ACK mechanism. Finally, QUIC is not fair with TCP, consuming more bottleneck bandwidth than the fair share, and has performance issues in older hardware, given that packet processing and cryptographic operations are not as optimized in the user-space as they are in the kernel. An implementation of QUIC for the open source network simulator ns-3 has also been recently released [75].

#### B. Stream Control Transmission Protocol (SCTP)

SCTP, which has been defined in [12], is a standardized signaling protocol developed by the IETF Signaling Transport working group that later evolved into a general-purpose protocol at the transport layer. It is now maintained within IETF's Transport Area (TSVWG) working group [76]. Like TCP, it is implemented in the kernel of the operating system, and offers a reliable, point-to-point, and connection-oriented data transport service over IP networks [77]. SCTP was specifically designed to be TCP-friendly, and shares its window-based congestion control, error detection and retransmission

<sup>&</sup>lt;sup>4</sup>Search Latency is defined as the delay from the time when a user enters a search term to the time all the result content is delivered to the client, including embedded content [70].

7

mechanisms. Nonetheless, SCTP also incorporates several new features [78], such as:

- **multihoming** while a TCP connection is established between two sockets (identified by an IP address and a port number each) at the two end hosts, an SCTP *association* can span across multiple IP addresses and possibly comprise several types of links in the network. For each client, a primary address is used to exchange data, and the others serve as backup in case of link failure or network fault, in order to ensure higher availability without interrupting the ongoing data transfers. Periodic *heartbeat* packets are sent to the backup addresses to verify the link state, and the retransmission of lost packets can also occur using backup addresses.
- **multistreaming** like QUIC, SCTP introduces multistreaming capabilities. Using the terminology from [12], it is possible to define an SCTP stream as a unidirectional logical data flow within an SCTP association [79]. SCTP then allows to split the application layer data into multiple substreams, according to the specific application that generates the packets. The sequencing and delivery are performed independently within each substream, thus overcoming the HoL blocking issue, previously described in Sec. II-C. Moreover, the delivery within each stream is strictly in order by default, as in TCP, but it is possible to enable unordered delivery [77].

The connection establishment and shutdown slightly differ from TCP to improve operations and security. A cookie mechanism during the initial 4-way handshake prevents SYN flooding<sup>5</sup> and masquerade attacks; moreover half-closed connections where one endpoint is still allowed to send data are not allowed.

The performance of SCTP has been evaluated in several scenarios: [80] reports that partial order and partial reliability increase the performance over traditional transfer protocols such as TCP. Jungmaier *et al.* [81] proved that SCTP traffic has the same impact of standard TCP traffic, that is, the two protocols can coexist in the same network. Moreover, the multihoming scheme can achieve better throughput and faster network accident recovery when applied to wireless LANs [82].

Finally, some papers have tested the performance of SCTP with respect to HoL blocking issues. Evaluations based on the Session Initiation Protocol (SIP) protocol at the application layer, and SCTP, TCP and UDP as transport, showed how the HoL blocking avoidance mechanism of SCTP does not provide better performance in normal scenarios of concurrent traffic on the network, and only slightly decreases the end-to-end latency with a high level of packet losses [83]. In addition, when the Maximum Transmission Unit (MTU) value suddenly decreases, SCTP is unable to avoid IP fragmentation and NATs may not be able to correctly route fragmented datagrams.

<sup>5</sup>SYN flooding is a denial-of-service attack where an attacker sends a succession of TCP SYN requests to a server but never replies to SYN-ACK responses. These half-open connections allocate memory that is freed only later, but in the meantime the system is not able to reply to legitimate queries.

## C. Datagram Congestion Control Protocol (DCCP)

DCCP, introduced in [23], [53], [54], is a protocol that augments UDP with connection-oriented functionalities such as congestion control, while maintaining its message-oriented structure. Like normal UDP, it does not guarantee reliable delivery, and packets can be out of order. Therefore, it presents a viable alternative as transport for applications which do not need reliability, but at the same time do not wish to implement their own congestion control mechanism [84]. The main reason why this protocol was introduced was the rapid growth of latency-oriented applications relying on UDP connections, like online games or IP telephony. These real-time applications cannot afford the additional latency introduced by reliability mechanisms, but still need to perform congestion control to conform with [85] and avoid choking TCP flows when sharing a bottleneck with them [86]. DCCP provides a common framework that can perform congestion control for these applications, but also allows new mechanisms to be implemented in a straightforward manner.

DCCP establishes a bidirectional connection, which is logically composed by two half-duplex unidirectional connections [23]. Given that it is an unreliable transport protocol, DCCP does not perform retransmissions. Nonetheless, it needs to detect losses to perform congestion control [84]. Therefore, the protocol does not feature cumulative ACKs, but a perpacket sequence number which allows the detection of missing datagrams. The ACKs are also acknowledged, using the same mechanism, and this makes it possible to perform congestion control on these feedback packets as well. The features of the connection (such as the security mechanisms, or congestion control algorithms), must be negotiated when establishing the connection, using control fields embedded in the DCCP headers [23]. Therefore, it is possible to tune the congestion control in each single end-to-end connection. DCCP also natively provides ECN [87], thus supporting the advertisement of network congestion without necessarily recording packet losses.

## D. Open Challenges and Research Directions

The protocols presented in this section have been proposed in the attempt to solve a number of open issues and challenges related to TCP. Nonetheless, there still exists a main problem related to their wide-scale adoption in the Internet, which, as we discussed, is throttled by the ossification of the networking equipment [16]. QUIC has adopted a promising direction, thanks to a user-space implementation, and to the support of major web content providers, such as Google.

Another important trend in QUIC and SCTP is the support for updates at the IP layer and/or multi-homing, so that the end-to-end connection does not need to be broken whenever one of the endpoints changes its IP address. This is a useful feature in highly mobile wireless networks, where the user may frequently roam across different networks. Additionally, these protocols natively support or can be extended to support multiple paths, as we will discuss in detail in Sec. V.

Moreover, QUIC, SCTP and DCCP share flexibility as a main design characteristic, which represents a novelty with This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/COMST.2019.2932905, IEEE Communications Surveys & Tutorials

8

Туре	Algorithm	Congestion control mechanism	Pros	Cons
Loss-based	NewReno [88] BIC [89] CUBIC [90] Wave [91]	AIMD Binary search increase function Cubic CWND function Burst-based adaptation	Proven convergence, fairness Higher efficiency RTT-independent fairness Fairness and efficiency	Inefficient in LFNs Too aggressive High number of retransmissions Highly volatile RTT
Delay-based	Vegas [92], [93] Verus [94] Nimbus [95] LEDBAT [96]	RTT changes as congestion signal Delay profile-based AIMD Explicit queue and cross traffic modeling Extra delay to high-priority flows	Fewer retransmissions, lower latency Adaptability to volatile channels Scalable aggressiveness to deal with CUBIC Does not affect other flows	Suppressed by loss-based flows High sender-side CPU load Unfair to Vegas and BBR Limited to low-priority traffic
Capacity-based	Westwood [97] Sprout [98]	Bandwidth estimation to decrease CWND HMM capacity model	Good in wireless and lossy links Low delay, customizable	No latency control Needs one buffer per flow
Hybrid	Compound [99] Illinois [100] Veno [101] BBR [21]	Sum of Reno and Vegas windows Delay used to determine CWND change Explicit model of the buffer Capacity and RTT measurement	Fast in LFNs, fair to CUBIC Fast in LFNs, fair to CUBIC Fast in LFNs, fair to CUBIC High throughput, low delay	No latency control No latency control No latency control Fairness and mobility issues
Learning-based	Remy [102] TAO [103] PCC [104] TCP-RL [105] QTCP [106]	Monte Carlo-based policy Advancement on Remy Online experiments to determine CWND Reinforcement learning to determine CC algorithm Reinforcement learning to select CWND	Reaches capacity with low delay Fairness issues solved Good in high RTT networks Self-organizing capabilities Higher throughput than NewReno	Fairness issues with other TCPs Requires knowledge of the network Untested with bufferbloat Untested in highly dynamic environments Limited performance evaluation

Table II: Summary of the main presented congestion control schemes, divided by design philosophy

respect to TCP. For example, in QUIC and SCTP it is possible to disable reliable transmissions for certain streams, and QUIC and DCCP have mechanisms for feature and parameter negotiation during the connection handshake. This reconfigurability should be however matched by a proper evolution of the APIs to the application layer, so that developers can benefit from it [107]. Explicit QoS signaling represents a step further in this direction: in order to fully support new applications and services, endpoints will need to be able to specify their QoS requirements [108], as will QoS-aware schemes exploiting tools such as SDN [109].

#### IV. CONGESTION CONTROL ALGORITHMS

Congestion control is an essential function in modern networks [88], since the high volumes of traffic that needs to be delivered reliably can only be supported if senders limit their rate before flooding their connections and affecting other flows. It can be performed directly by transport layer protocols such as DCCP, SCTP, QUIC, and the omnipresent TCP. Applications that run directly over UDP usually perform their own congestion control, as recommended by the IETF [85].

The rationale behind congestion control is the following: since links with limited capacity need to be shared by multiple flows, they can get overloaded by the aggregate flow generated by multiple (or even single) sources, becoming the bottleneck of the connection. If no limit is set to the sending rate, senders will start retransmitting packets when they do not get a positive reception acknowledgment (ACK) from the destination within an interval called Retransmission Timeout (RTO). This, in turn, increases the load on the network, triggering a destructive spiral known as *congestion collapse* [110]. Congestion control is aimed at avoiding this condition: senders react to signs of congestion by backing off and reducing their own rate proactively. There are several congestion control design philosophies, with different ways to detect congestion and react to it, but they all share this common core.

The issues we mentioned in Sec. II make the design of good congestion control mechanisms extremely difficult; specific features of the connection and of each individual link can significantly affect the performance, often in unpredictable ways [34]. Mathematical models of protocols and their rigorous assumptions are rarely realistic, and often theoretically optimal algorithms do not work as expected in real networks [111]. Moreover, in most cases (e.g., for TCP), the congestion control algorithm for a certain transport protocol is implemented in the same kernel code base of that protocol and, consequently, it is the same for every end-to-end connection. Therefore, it is not possible to customize the response of the congestion algorithm to the characteristics of each connection.

Furthermore, detecting when congestion is taking place and how the congestion window size should be maintained are not trivial problems. Most early solutions adopted packet loss as a sign of congestion and used an AIMD strategy, which ensured asymptotic fairness and stability [112]. The AIMD scheme is simple, as it only needs two commands to update the congestion window (often referred to as CWND for short):

$$\text{CWND} \leftarrow \text{CWND} + \frac{\alpha}{\text{CWND}} \tag{1}$$

on received ACKs and

$$\text{CWND} \leftarrow \frac{\text{CWND}}{\beta} \tag{2}$$

on packet losses, where  $\alpha$  and  $\beta$  are mechanism-set parameters. Newer protocols often include delay information as a potential congestion signal, and their adaptation of the congestion window is more tailored to their operating requirements. Selective ACK (SACK) is another option that improves the information available to the congestion control

mechanism [55]. We now list the main congestion control design philosophies and their advantages and shortcomings, quickly reviewing older algorithms in order to provide the necessary background information to better appreciate the new developments; the main mechanisms presented in this section are summarized in Table II. We also highlight that all the TCP congestion control mechanisms discussed below can be implemented for OUIC and SCTP senders with no modifications. DCCP requires more effort, since it does not include a retransmission mechanism, but the basic principles and design philosophies are the same. DCCP provides a framework for congestion control [84] mechanisms which can be selected by the application through the Congestion Control IDs (CCIDs). The congestion control mechanisms are separated from the core of the protocol enabling modularity and possible expansions. Each one of the two half-connections can even choose a different CCID since they are logically separated. In the following, we will refer to the standardized DCCP congestion control schemes by their CCID. Another taxonomy of congestion control mechanisms for TCP was presented in [113], but lacks the most recent developments (e.g., those related to machine-learning-based algorithms) that we survey in this paper.

## A. Loss-based mechanisms

Classic TCP congestion control algorithms such as Tahoe [1], Reno [114] and New Reno [88] use packet losses to detect congestion. If a packet is determined to be lost, either when the RTO timer expires or after three consecutive duplicate ACKs, the AIMD mechanism sharply decreases the congestion window. While Tahoe is very conservative, starting back from a CWND of 1 packet and entering the Slow Start phase, Reno and New Reno only reduce CWND by half in case of three duplicate ACKs, as shown in Fig. 8. The default QUIC congestion control is also based on New Reno [88], even though other congestion control algorithms may be used, as in TCP [48]. A number of refinements that have been proposed for New Reno are directly included in the IETF QUIC protocol drafts. For example, the packet number is monotonically increasing, so that the QUIC sender knows if a received ACK is for the original packet or its retransmission. Moreover, while with the TCP SACK feature only three ACK ranges can be specified [115], QUIC does not have a limit for the number of ACK ranges in a feedback packet. The draft [48] also recommends the use of pacing, to avoid the bursty transmission of packets, and a tail loss probe mechanism to efficiently detect packet loss at the end of a data transfer.

The BDP is the product of the bottleneck capacity and the minimum RTT; networks with a BDP over 100 kb are considered Long Fat Networks (LFNs) [116]. The classic lossbased mechanisms are extremely inefficient in LFNs, since the recovery phase will take slightly less than a minute for a connection with a minimum RTT of 100 ms and a capacity of 100 Mb/s, and about 10 minutes if the capacity is increased to 1 Gb/s [89].

Binary Increase Control (BIC) [89] and CUBIC [90] are two congestion control algorithms developed to solve earlier



9

Figure 8: Evolution of TCP New Reno's CWND over time

mechanisms' issues with LFNs. They abandon pure AIMD for a more complex function (a binary search in BIC's case, and a cubic function of the time since the last loss for CUBIC, as shown in Fig. 9) in order to achieve fairness between flows with different RTTs. However, BIC was deemed too aggressive and unfair to legacy flows [117], while CUBIC is now widely deployed [118] (it is the default congestion control algorithm in the Linux kernel since version 2.6.19) and is currently an IETF Internet-Draft [119] because it can achieve high performance without being unfair to flows using older congestion control mechanisms. The QUIC implementation in the Chromium code base also features CUBIC congestion control [70], [90]. The main difference between CUBIC applied to TCP and QUIC is that  $\beta_{\text{TCP}} = 0.3$  (i.e., the decrease factor for the congestion window after loss events), while  $\beta_{\text{OUIC}} = \beta_{\text{TCP}}/2$  [72].

DCCP implements three loss-based mechanisms: CCID 2 [120] is a New Reno-like mechanism that uses the same AIMD principle, but using packets as a data unit instead of bytes because of the datagram-oriented nature of the underlying UDP socket. CCID 2 also applies congestion control to acknowledgments, and has no retransmission mechanism since it is meant for applications that do not require reliable delivery and might even be hampered by it. CCID 3, or TCP-Friendly Rate Control (TFRC), uses the TCP throughput equation explicitly [121] to calculate the sending rate in a way that is both fair to competing TCP flows and relatively stable. It is recommended for applications that need to maintain a more stable throughput while remaining TCP-friendly. CCID 4 is a variant of TFRC designed for applications that send small packets, achieving roughly the same bandwidth as an equivalent TCP flows sending full-sized packets [122]. Another extension of TFRC, which behaves better when multiple DCCP flows share a bottleneck, was proposed in [123].

Nowadays, most of the research community is moving away from loss-based congestion estimation. Several simulation studies show that none of these mechanisms performs well in wireless networks [118], in Mobile Ad Hoc Networks (MANETs) [124] and in broadband, high-latency networks, [125]. TCP, QUIC, and SCTP all suffer from the same issues, since the problem is inherent in the loss-based congestion control philosophy. However, it can still be a useful



Figure 9: CWND growth of CUBIC algorithm

design principles for some network scenarios, particularly when very high RTTs are involved. One example is TCP Noordwijk [126], a cross-layer version of TCP that exploits the peculiarities of satellite links. The main feature of TCP Noordwijk is that it transmits data in bursts at the maximum capacity of the bottleneck, assuming a large buffer at the bottleneck (as is the case for satellite links) and not considering competing flows. Fairness is ensured by appropriately scheduling transmission bursts; a recent proposal, called TCP Wave [91], extends the protocol by removing the cross-layer aspects and adapting the algorithm to any kind of link. The protocol quickly achieves fairness with any other loss-based flow, and goes back to full bandwidth utilization just as quickly if the competing flows stop. However, if the protocol sends packets in bursts, packets will be queued at the sender until the next burst, causing a longer queuing time. Burstbased mechanisms are also highly vulnerable to jitter and capacity fluctuations because they need to commit to a sending rate in advance and have fewer opportunities to correct their mistakes than schemes which make smaller commitments: if the bottleneck link is time-varying, they take longer to react and can make bigger errorsMany common applications are sensitive to both delay and jitter, reducing the possible usefulness of Wave in standard networks.

#### B. Delay-based mechanisms

Loss-based mechanisms such as CUBIC can efficiently make use of most connections; however, the only signal of congestion they perceive is, as the name suggests, packet loss. If the buffer at the bottleneck link is large, packets will accumulate until the buffer fills, while the sender will keep increasing its sending rate and the queueing delay will keep rising. Latency can even reach 10 seconds in the conditions of heavy bufferbloat we talked about in Sec. II-A [68], and filling up the buffer completely can cause instability in the throughput observed by other flows.

One possible solution to this is delay-based congestion control: losses may be a signal of congestion, but congestion may occur long before the buffer is full, and its early detection can increase the reactivity of the congestion control, which can then back off as soon as the RTT increases substantially and avoid sharp declines in the throughput and high latency.

TCP Vegas [92] was the first implementation to use delay as a signal of congestion: after a milder slow start phase [127], it adjusts the congestion window to the expected throughput of the channel. The channel is assumed to be underused if the RTT is close to the minimum measured one; as soon as the RTT grows beyond a certain threshold, the protocol backs off. Vegas is fair and uses capacity better than Reno [92] while maintaining a low latency. However, the first version of Vegas had severe issues with the Delayed ACK mechanism, which affected its RTT estimation, and with links with high latency [112]. These issues prompted the development of several improvements, like Vegas-A [128], which addresses rerouting and bandwidth sharing fairness, Vegas-V [129], which increases the aggressiveness of the original algorithm, while remaining fair to other flows, and Adaptive Vegas [130], where *adaptive* refers to the ability of the algorithm to change its parameters according to the statistics of throughput and RTT variation. Some specialized versions were developed for wireless [131] and mobile [132] ad hoc networks; three particular modifications became part of the so called New Vegas algorithm [133]. New Vegas [93] uses packet pacing in the slow start phase to avoid issues caused by burstiness, and sends packets in pairs to avoid the Delayed ACK problem, while maintaining most of Vegas's main ideas.

A simple delay-based congestion control mechanism for DCCP has been proposed in [134], reducing the sending rate by the ratio between the minimum and current RTT. A similar Vegas-like mechanism [135] uses one-way delay to measure congestion, adapting the send rate to maintain QoS. Since these schemes have not been standardized, they were not assigned a specific CCID.

Over the last few years, the development of delay-based protocols has seen a new renaissance thanks to the rise of interactive applications with strict latency requirements. Verus [94] is a recent delay-based end-to-end congestion control protocol that exploits delay measurements to quickly adapt the congestion window size to enhance the transport performance in cellular networks. It retains the same AIMD scheme as traditional TCP, but changes the additive increase part: by constantly sensing the channel, the window size is increased or reduced at every step according to the observation of shortterm packet delay variations, that the authors called "learning a delay profile." Basically, the protocol builds an empirical function by associating CWND and delay and observing when the rate exceeds the capacity of the channel. The multiplicative decrease and slow start steps remain unmodified. This strategy can be more effective than using an explicit model of the connection, and is practically achieved by using a sliding window over a period equal to the estimated RTT.

Real world tests were performed by the authors of [94] on 3G and LTE networks considering multiple competing flows between some devices and a server on a high bandwidth and low delay network, whose results show that Verus can achieve a throughput similar to that of TCP CUBIC while reducing the delay by an order of magnitude. The authors also pointed out how Verus outperforms other protocols and quickly adapts to rapid network changes and to the arrival of new flows, with throughput being independent of the RTT, confirming the protocol's fairness. While the positive results were confirmed by independent measurements [30], these tests also highlighted the massive computational load that the Verus delay profile estimation imposes on the sender; its use in uplink or highthroughput scenarios might be impossible because of the protocol's complexity.

In general, delay-based mechanisms are more stable, retransmit fewer packets and have lower latencies than lossbased mechanisms, with a similar throughput in realistic conditions; however, their adoption has been blocked because of the higher aggressiveness of the existing protocols. If a delay-based and a loss-based flow share a bottleneck, and the buffer at the bottleneck node is large enough, the delay-based mechanism will sense congestion far before its loss-based competitor, which will keep increasing its sending rate [136]. As a result, the throughput of the delay-based flow will basically drop to zero, while the loss-based flow will take up all the available capacity. Since the majority of the servers in the Internet run a loss-based congestion control mechanism, and separating loss-based and delay-based traffic was extremely expensive and required extensive changes, delay-based protocols have never been deployed on a wide scale. However, deploying better performing delay-based protocols, that might suffer from the aggressiveness of loss-based flows, has recently become possible thanks to techniques such as network slicing and Network Function Virtualization (NFV) [137], although it still involves a significant effort on the network operator's part.

Two delay-based protocols that can switch to a more aggressive mode to avoid being outcompeted by loss-based flows have recently been proposed. Like Verus, Copa [138] is a protocol that uses the delay profile to determine the sending rate. It uses Markov chains to explicitly model the bottleneck queue, and it dynamically adjusts its aggressiveness to compete fairly with loss-based protocols. Nimbus [95] is another enhancement that also models the elasticity of cross traffic, detecting how other flows will react to changes in the perceived bandwidth and adapting to the resulting scenario. The cross traffic modeling is performed by observing the Fourier transform of the capacity and detecting periodic behavior; this approach works well with CUBIC or Copa cross-traffic, but can fail for BBR cross-traffic and is too aggressive in the presence of Vegas flows, which it senses as elastic flows and proceeds to outcompete.

Another interesting development is represented by TCP Low Extra Delay Background Transport (LEDBAT) [96], a delaybased algorithm developed for BitTorrent traffic: since background traffic should have a lower priority than user traffic, the low aggressiveness of delay-based mechanisms becomes a strength. LEDBAT estimates the capacity left unused by foreground flows by measuring the extra delay it adds after sending a packet, and it only transmits if it deems its actions will not affect more important flows. It is currently deployed on Microsoft Windows machines as a solution for bulk transfer applications. TIMELY [139] is another delay-based protocol which adjusts the congestion window using the gradient of the RTT, designed specifically for data centers. Finally, TCP LoLa [22] adapts the CUBIC mechanism to a delay-based detection method; when the queueing delay starts to increase, LoLa switches to a Vegas-like holding mechanism to maintain capacity and fairness to other flows.

## C. Capacity-based mechanisms

TCP Westwood [97] was a congestion control mechanism from 2001 that explicitly tracked the estimated bandwidth in order to adaptively change the congestion window after a loss, while maintaining the additive increase scheme of TCP Reno. Westwood was designed for wireless and lossy links [140], which can have a fast-varying capacity, as well as packet losses caused by the physical layer. Both of these factors can affect transport layer performance, as traditional loss-based schemes will overreact to losses by halving the send rate even when there is no need for it, and fast-varying capacity can also be interpreted as a sign of congestion. Westwood adapts the congestion window to the measured connection capacity after losses and achieves a higher performance [141] than loss-based schemes in those conditions, as well as in wired networks.

However, the idea of capacity-based congestion control was recently revived by Sprout [98], an end-to-end protocol designed for cellular wireless networks that aims to be a compromise between achieving the highest possible throughput and preventing long packet delays in a network queue. Sprout exploits some of the peculiarities of cellular networks, improving congestion control in that particular scenario. Since cellular networks often suffer from extreme bufferbloat, using packet loss to sense congestion leads to extremely high latency. Additionally, cellular network users are not subject to queues accumulated by other users, since carriers generally set up separate uplink and downlink queues for each device belonging to a cell. Sprout exploits that fact by periodically measuring the connection capacity and predicting its future distribution with a Hidden Markov Model (HMM). Estimating the available capacity is instead done by counting the received bytes in a long interval and then dividing the result by its duration. This estimate is used to forecast the number of packets that is safe to send over the links, that is, the right number of packets so that they will not wait too long in a queue. It is a highly customizable protocol: the model can be more throughput or delay-oriented by changing an aggressiveness parameter. However, this approach requires a dedicated buffer: if multiple flows share the same buffer, Sprout will leave almost all the capacity to the competing flows [30].

## D. Hybrid mechanisms

Some congestion control mechanisms try to combine two approaches to reap the benefits of both. Compound TCP [99] is a well-known example, as it is available on all Microsoft Windows machines by default. Compound uses the sum of a delay-based window and a loss-based Reno window as its CWND, and as such it can deal with LFNs better than pure loss-based mechanisms while maintaining a measure of fairness toward them, since its compound window will grow faster than a purely loss-based one, but the protocol's behavior will revert to loss-based congestion control when it senses congestion. Illinois [100] is another algorithm that considers



Figure 10: BBR basic principle: the mechanism tries to work at the optimal operating point instead of filling the bottleneck buffer

both loss and delay to infer congestion: while it uses losses to detect it, the variation in the delay is used to determine the rate of change in the congestion window. It is friendly to CUBIC, but, like Compound, its CWND grows significantly faster when it is far lower than the capacity; Veno [101], a hybrid between Vegas and Reno, uses an explicit model of the bottleneck buffer occupancy to achieve the same kind of CWND growth.

The most notable recent development in hybrid congestion control schemes is undoubtedly Google's BBR algorithm, which takes a different approach by trying to estimate both the bandwidth and the RTT [21]. It tries to reach the optimal operating point [142] by keeping a CWND equal to the BDP; any further increase in the CWND just increases the delay without any throughput benefits, as shown in Fig. 10. While converging to the BDP with traditional mechanisms is impossible [143], BBR alternately estimates the minimum RTT and the capacity by periodically draining the buffer, then measuring the RTT, and finally returning to the initial delivery rate. This mechanism is shown in Fig. 11.

BBR can share the bandwidth fairly with other BBR or loss-based flows [21], but its latency is low only when the connection is used by other delay-conscious protocols, as there is no way to seize a fair share of the capacity from lossbased flows without having significant queueing delays. BBR also works well with high-capacity connections; however, the presence of short buffers causes the BBR operating point to increase too much, thus leading to massive packet losses and unfairness towards loss-based flows [144].

Fairness and stability when sharing a connection with other congestion control algorithms are two of BBR's biggest open issues. Competition with other algorithms can cause severe throughput fluctuations [145], and BBR's aggressiveness against loss-based flows strongly depends on the buffer size: as previously mentioned, the mechanism becomes too aggressive when buffers are very small [144] and too conservative when the network suffers from bufferbloat. The same problem was noticed by Farrow in his comparative study of various congestion control algorithms, namely CUBIC, NewReno and BBR, in heterogeneous network environments [146]: virtual tests with several competing flows both with the same and with different algorithms lead the author to conclude that NewReno and CUBIC are sufficiently fair, with the latter offering more predictable and stable throughput, while asserting that BBR is unfair and still not ready for public use.

Doubts have also been raised about BBR's performance in cellular networks [147] and under mobility [148], [149], but work to adapt it to this kind of environment is currently ongoing. BBR has also been proposed for QUIC [150].

#### E. Machine learning approaches

Over the past few years, machine learning has become an important tool for network and protocol designers [151]; research on machine-learning based congestion control is ongoing, and there already are several working mechanisms in the literature.

TCP Remy [102] is the first example of machine learningbased congestion control: the authors define a Markov model of the channel and an objective function with a parameter  $\alpha$ which can be tuned to set the aggressiveness of the protocol  $(\alpha = 1 \text{ corresponds to proportional fairness, while } \alpha = 0$ does not consider fairness at all and  $\alpha = \infty$  achieves maxmin fairness), and use a machine learning algorithm to define the behavior of the congestion control mechanism. The inputs given to the mechanism are the ratio between the most recent RTT and the lowest measured RTT during the connection, an Exponentially Weighted Moving Average (EWMA) of the interarrival times of the latest ACKs, and an EWMA of the sending times of those same packets. The mechanism itself is essentially a Monte Carlo simulation, and has some limitations, since it requires some prior knowledge about the network scenario and operates on the basic assumption that all competing flows use it. A more advanced version of Remy, called Tractable Attempt at Optimal (TAO) [103], solves the problem of TCP awareness and performs well with heterogeneous competing flows, but still requires extensive prior knowledge about the network to function.

Performance-oriented Congestion Control (PCC) [104] is another learning-oriented mechanism that uses online experiments instead of offline pre-training: it uses SACKs to measure the utility of an action, and adjusts the sending rate accordingly. PCC is more aggressive than even loss-based TCP versions, and performs well with short buffers and high RTTs, but its performance in wireless networks or links affected by bufferbloat has not been tested. PCC Vivace [152] is a new version of the congestion control mechanism with improved TCP friendliness and a more advanced learning mechanism using linear regression techniques.

An attempt at using more advanced supervised learning techniques to design congestion control mechanisms was made in [153], and Q-learning was used in [154] and [155]. However, all these algorithms are optimized for very simple network topologies and situations and strongly depend on the considered scenario, so to the best of our knowledge there is no fully implemented and well-tested congestion control mechanism using these techniques. Along this line, the authors of [105] investigate a reinforcement learning approach to improve the performance of both short and long TCP flows.

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/COMST.2019.2932905, IEEE Communications Surveys & Tutorials



Figure 11: Overview of the operating points in the states of the BBR algorithm

In particular, the proposed approach, TCP-RL, exploits two different learning processes to adapt the initial window of the connection (i.e., the number of packets that can be sent when the connection starts) and the congestion control algorithm used for each single flow. The performance evaluation is based on a real implementation, and either real traffic or synthetic traces generated from a major web service provider. Finally, also QTCP [106] adopts a reinforcement learning approach, but the actions that are available are directly related to the dimensioning of the congestion window. The authors, however, compare the performance of their scheme with TCP NewReno.

A more detailed discussion on the design considerations for new machine learning-based congestion control mechanisms can be found in [156].

#### F. Cross-layer approaches

As previously mentioned, TCP relies on an abstraction of the multiple hops between the two end-points of the connection to perform its congestion control operations. This, however, is sometimes sub-optimal, given that TCP may fail to capture medium-specific characteristics (e.g., mmWave link variations, as discussed in Sec. II-D). Therefore, cross-layer approaches aim at designing congestion control algorithms that exploit additional and more precise information provided by the lower layers of the protocol stack, and are popular especially in the wireless domain.

Several papers have recently proposed cross-layer solutions to cope with cellular networks operating at mmWave frequencies. In [157], [158] the authors tune the TCP congestion window to the BDP of the connection, tracked with physical and Medium Access Control (MAC) layer information in the mobile terminal, for downlink and uplink TCP connections, respectively. The paper [159] proposes a proxy-based mechanism at the base station that intercepts ACKs and changes the advertised window value to force the unmodified TCP sender to track the BDP of the end-to-end connection with a mmWave link.

The SDN paradigm has also opened new possibilities for cross-layer approaches, giving the transport layer the possibility to reserve resources for low-latency applications [160] by communicating directly with the network controller.

Another possible cross-layer approach is to consider the needs of the application: several TCP versions that consider rate distortion and content type in multimedia streaming have been proposed [161]–[163].

In [164], the authors revisited a cross-layer approach for 3G and 4G cellular networks applying it to QUIC's congestion control, with modifications both at the sender (which alternates

between a rate-limited approach, when a valid rate estimation is available from the lower layers, and a congestion-windowlimited state) and at the receiver (which feeds back the rate estimation to the sender).

13

#### G. Datacenter networks and the Incast issue

As we described in Sec. II-B, data centers are a very peculiar networking environment, which requires several adjustments on the transport layer to avoid issues like the Incast problem. Over the past few years, there have been several attempts to mitigate these issues. DCTCP [45] is a congestion control algorithm based on ECN, which reduces the congestion window based on the amount of experienced congestion, leaving enough buffer space to avoid the Incast issue.

Another option is to abandon standard congestion control altogether and use explicit rate control, as in the Deadline-Driven Delivery (D<sup>3</sup>) scheme [165], which can solve the problem if full network support is available. This can be implemented by using SDN and NFV solutions [166], but scalability remains an issue, and D<sup>3</sup> is vulnerable to bursts. Deadlineaware Data center TCP (D<sup>2</sup>TCP) [167] is a distributed scheme that considers application-level deadlines for data blocks like D<sup>3</sup>, while avoiding its need for full network support and vulnerability to bursty traffic. It is similar to DCTCP, but it has full awareness of the deadlines when performing congestion avoidance.

The pFabric scheme [168] is another network support technique based on early dropping: switches implement priority dropping with very short buffers, enforcing flow scheduling even if the flows are aggressive and do not back off until the loss rate is consistently high. In this kind of scenario, flows need to be aggressive, or the early dropping mechanism would reduce their throughput far below the achievable limit. The pHost [169] protocol aims at replicating pFabric's performance without any network support, using a token-based scheme at the receivers to assign priorities to senders. Receivers prioritize the flow with fewest remaining bytes when assigning tokens, achieving near optimal performance. ExpressPass [170] is another technique that uses receiver-side credit packets to control sender-side rate; credit packets can be lost without consequence, and the loss rate is used to gauge the connection capacity. Homa [171] is a new connectionless protocol inspired by pHost and ExpressPass, which can significantly reduce latency with no network support. It aggressively uses priority queues on switches, using receiver-driven flow control to assign the priorities and explicitly limiting the Incast issue by counting outstanding requests and marking new packets to reduce their priority if they are over a certain length. Its connectionless nature limits HoL blocking and avoids explicit acknowledgments, further contributing to the protocol's efficiency.

Prioritization, Arbitration, and Self-adjusting Endpoints (PASE) [172] is a strategy that synthesizes the approaches of DCTCP,  $D^3$ , and pFabric. It uses all three mechanisms at different scales: explicit rate control is performed at coarse time-scales and with low precision, since the endpoints themselves can find an efficient allocation by probing according to their assigned priority. Finally, pFabric-like prioritization and dropping are performed at very short timescales. This hybrid approach can outperform each adaptation by itself.

For a more complete discussion of data center networks and their transport layer issues, we refer the reader to [33], [173]–[175].

## H. Open Challenges and Research Directions

Congestion control is a critical component in the operation of most transport protocols. TCP, SCTP, and QUIC implement it natively, and DCCP represents an evolution in the design of unreliable transport protocols that also accounts for congestion. The research in congestion control has seen a fast evolution in several directions: while traditional loss-, delay-, and capacity-based approaches have all been used as stepping stones towards more efficient protocols, entirely new machine learning-based and cross-layer protocols are beginning to attract significant interest.

The main trend in this area is related to overcoming a onesize-fits-all approach to congestion control (e.g., the widely adopted NewReno and CUBIC) towards algorithms that target a particular use case. For example, several delay-based protocols target low queuing delay (e.g., TCP Copa, LoLa), which translates into low end-to-end latency. However, updates in the TCP congestion control mechanism require modifications in the kernel, and thus an Operating System update, which is a cumbersome operation. The development of user-space protocols, such as QUIC, allows researchers to experiment more, and network support through SDN and slicing makes it easier to develop mechanisms that are not as aggressive as CUBIC, since they can be isolated in separate buffers and thus do not need to compete with it.

#### V. MULTIPATH TRANSPORT PROTOCOLS

Nowadays, most devices can use multiple communication technologies at the same time: for example, modern smartphones can connect to both Wi-Fi and LTE. For this reason, multipath communications have become the subject of considerable interest over the past few years. At the transport layer, multipath-capable protocols need to be designed to successfully exploit the advantages of multipath diversity, but this is not always simple.

## A. MPTCP

MPTCP is a fully backward-compatible extension of TCP, published as an experimental standard in 2011 [25], [26], [176] and now widely deployed [177]. It allows applications to use



Figure 12: Protocol stack with MPTCP.

multiple connections at the same time without any changes to the socket API. Using multiple connections can improve the total capacity, provide redundancy against link failures, and reduce the load on congested paths. In order to ensure backward compatibility and transparency to applications [25], MPTCP needs to be implemented within the operating system's kernel.

Fig. 12 shows the basic architecture of an MPTCP host: the connection is composed of two separate TCP flows on different paths [178], each with its own congestion control and ACKing mechanism. Using single-path TCP flows is necessary because many middleboxes inspect TCP traffic and discard packets with behaviors inconsistent with the normal operation of the protocol (e.g., missing ACKs, out-of-order or withgaps sequence numbers, or badly formed options) for security reasons.

An MPTCP session starts with a single TCP sub-flow; MPTCP-capable hosts can recognize each other by setting the MP\_CAPABLE option in the first packets of the connection handshake. The two hosts then exchange cryptographic keys to add new sub-flows securely; these can be established by setting the MP\_JOIN option and using a hash of the connection's keys during the standard TCP handshake. MPTCP supports also the addition and removal of addresses on a host, both implicitly and explicitly.

MPTCP assigns packets two sequence numbers: the standard sequence number for each sub-flow and a connectionlevel Data Sequence Number (DSN). Since congestion control and retransmission need to happen at the sub-flow level, as each sub-flow is a full-fledged TCP connection in its own right, the flow-level sequence numbers need to be consecutive and data on each sub-flow is delivered in-order; the connectionlevel DSN is needed to join the two split data streams and reconstruct the original data [179]. Retransmissions of the same data over multiple sub-flows are also possible, since the connection-level DSN can identify the packet at the connection level and mitigate HoL if a sub-flow is blocked. However, data transmitted on a sub-flow must also be retransmitted and

![](_page_14_Figure_2.jpeg)

Figure 13: Scenarios with unfairness between MPTCP (with uncoupled congestion control) and TCP, adapted from [32].

delivered on it, even if received correctly on another, since doing otherwise would break the TCP-compliance of the subflow.

As with sequence numbers, there are two types of ACKs: regular acknowledgments on each sub-flow, and connectionlevel ACKs, which act as a cumulative acknowledgment for the whole data flow.

The support of multiple IP addresses for the same MPTCP connection may result in security vulnerabilities [180], but the design of the protocol minimizes this risk, especially if we compare it to single-path TCP, which is vulnerable for example to the man-in-the-middle attack as well. Multiple IP support can also be a problem for applications, since the application cannot implicitly rely on the one-to-one mapping between host and address and has to adapt its code to deal with multiple IP addresses [179], [181].

MPTCP's main rationale is that using multiple flows will reduce delay and increase throughput and reliability; the protocol is now implemented in most major operating systems and used in data centers and wireless networks [182]. However, several issues need to be taken into consideration, and the multipath scheduler and congestion control mechanism need to be carefully tuned, particularly in wireless networks [67], [183].

The HoL problem we described in Sec. II-C is particularly severe in MPTCP: since each sub-flow forwards data to the multipath level only in-order, a loss on one sub-flow can block the whole connection until the packet is retransmitted on that same sub-flow. If the loss happens on the sub-flow with the longest RTT, this can take a long time, particularly in networks affected by bufferbloat. HoL can even reduce the throughput of other sub-flows, since the receiver-side buffer is limited and once it is filled other flows will have to wait for the head of line packet without transmitting anything [184]. In some cases, MPTCP is clearly outperformed by a single-path TCP flow on the best path [67]; it has even been proposed to completely disable it and automatically fall back to single-path transmission in lossy scenarios [185].

1) Congestion control: The first versions of MPTCP used *uncoupled* congestion control: each path had an independent

congestion window, updated with one of the single-path mechanisms we described in Sec. IV. However, measurements show that uncoupled MPTCP can be unfair to single-path TCP users sharing one of the paths [186], and can often lead to increased resource usage without a corresponding performance benefit [189].

In general, a congestion control algorithm for multipath flows should [26]:

- 1) Obtain at least the same throughput as a single path connection on the best sub-path;
- Not take up more capacity than what a single path on the best sub-path would;
- 3) Avoid to push traffic on the most congested sub-paths.

Fig. 13 shows two different scenarios which lead to unfairness when MPTCP applies uncoupled congestion control [32]: in the first one, the two paths of the MPTCP share a bottleneck. In this case, goal 2 is violated, since MPTCP behaves like two separate TCP flows, taking up two thirds of capacity and squeezing out the competing single-path flow. In the second scenario, one of the two paths shares its bottleneck with a standard TCP flow. In this case, MPTCP violates goals 2 and 3: since it has a non-congested path, MPTCP should steer most of its traffic towards the free path, avoiding the congested one (goal 3) and taking up a total capacity of R (which would respect goal 2).

In order to achieve fairness with single-path flows, a *semicoupled* algorithm achieving good performance while being fair to single-path flows, called Linked Increases Algorithm (LIA), was proposed [188] and subsequently standardized [26]; the fully *coupled* version [187] was considered too conservative, as it did not use the most congested path at all and could not then discover any changes in the traffic on that path.

The LIA congestion mechanism maintains the uncoupled single-path slow start, fast recovery and fast retransmit mechanisms of single-path TCP, but increases the congestion window after each ACK received on sub-flow i by

$$\min\left(\frac{\alpha \times \text{bytes}_{\text{acked}} \times MSS_i}{CWND_{\text{total}}}, \frac{\text{bytes}_{\text{acked}} \times MSS_i}{CWND_i}\right), \quad (3)$$

where the second argument in the minimum is the increase that an uncoupled TCP flow would have, while the first one is a coupled value that takes the overall congestion window into account. The parameter  $\alpha$  describes the aggressiveness of the algorithm. Theoretically, we get

$$\alpha = CWND_{\text{total}} \ \frac{\max_i (CWND_i/RTT_i^2)}{\sum_i (CWND_i/RTT_i)^2}, \tag{4}$$

which corresponds to the congestion window increase that satisfies goal 1: by setting that value of  $\alpha$ , the MPTCP flow can theoretically obtain the same throughput as a single-path flow on the best path. The complexity of calculating the congestion window or  $\alpha$  for each ACK scales linearly with the number of flows.

The Opportunistic Linked Increases Algorithm (OLIA) [189] is a modification to the LIA algorithm that achieves optimal resource pooling and avoids LIA's fairness issues. Balanced Link Adaptation (BALIA) [190] is

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/COMST.2019.2932905, IEEE Communications Surveys & Tutorials

Table III: Summary of the	main presented congestion	control schemes for multipath,	divided by protocol
---------------------------	---------------------------	--------------------------------	---------------------

Protocol	Algorithm	Mechanism	Туре	Pros	Cons
МРТСР	CUBIC [186] Coupled [187] LIA [26], [188] OLIA [189] BALIA [190] BELIA [191]	Independent CUBIC Joint CUBIC CWND Weighted CWND increase Resource pooling Mix of LIA and OLIA Capacity-based	Uncoupled Fully coupled Semi-coupled Semi-coupled Semi-coupled	Direct adaptation of TCP Fair to TCP Fair to TCP, uses all paths Pareto optimal Mix of OLIA and LIA's pros Good in wireless networks	Unfair to TCP, HoL Does not use most congested path HoL, unfairness to other LIA flows HoL, unstable behavior HoL, wireless networks HoL, unfairness
	DRL-CC [192]	Reinforcement learning	Fully coupled	High fairness	Untested in wireless networks
SCTP	CMT [193] MPTCP-like [194]	SACK-based CWND growth Adaptation of LIA	Semi-coupled Semi-coupled	No HoL No HoL, fair to TCP	Retransmissions, wireless networks Retransmissions, unfairness
MP-DCCP	FSCC [195]	Adaptation of CCID 2	Semi-coupled	Free paths help congested ones	Untested, high reordering
LEAP	CKF [30]	Capacity-based Kalman filtering	Uncoupled	No HoL, good in wireless networks	Vegas-like conservativeness

another proposal, now widely adopted, that tries to balance TCP-friendliness and throughput stability, striking a balance between LIA's and OLIA's strengths and weaknesses.

There are also proposals for delay-based multipath congestion control, either by using delay-based algorithms on each sub-flow or by using specifically designed solutions [206]. The same goes for capacity-based [191], [207] and cross-layer schemes [208], but all these proposals are untested and still under development.

However, these algorithms have been shown to be inefficient in wireless networks [209], [210]: the HoL problem is exacerbated by the volatility of wireless links, and errors and retransmissions negatively affect both the average throughput and its stability, resulting in an oscillatory behavior [211]. A promising approach to multipath congestion control is to use reinforcement learning tools: Deep Reinforcement Learning Congestion Control (DRL-CC) [192], an application of the actor-critic method which jointly sets the congestion window for all active flows and all paths, achieves high fairness in a wired network scenario with multiple active flows. The performance of these types of algorithms in more volatile wireless environments is still untested.

Forward Error Correction (FEC) has been proposed as a solution for the HoL problem: by adding some redundancy packets to the flow, errors can be recovered by the receiver by decoding them and avoiding blocking other sub-flows unnecessarily. However, each sub-flow still needs to retransmit lost packets and deliver everything in order, even if the data has already been received on other flows. MPTCP with Systematic Coding (SC-MPTCP) is a hybrid solution [212] that uses FEC to reconstruct missing packets, and other recent solutions [184], [213]-[216] use different coding schemes to achieve the same objective. These schemes can mitigate the HoL problem, especially when the receiver buffer size is limited, but they often require specifically designed schedulers and do not fix other congestion control issues in multipath, such as throughput volatility [217]. Another workaround for the HoL problem is to retransmit lost packets on a faster path while avoiding transmitting new packets on the congested path [218].

The main multipath congestion control algorithms we presented are summarized in Table III; for a more thorough survey of the research on multipath congestion control and MPTCP, we refer the reader to [32]. 2) *Scheduling:* Congestion control is not the only factor affecting performance, as scheduling is also extremely important in MPTCP: sending data on the wrong path can exacerbate the HoL problem, increase latency and lower throughput.

The scheduler in the Linux kernel currently follows the Lowest-RTT-First (LowRTT) policy, so the first packet in the send buffer will always be sent through the lowest RTT available path, but this simple heuristic is not always efficient. As Hwang *et al.* point out in [219], waiting until the fastest path becomes free could be more efficient than immediately sending a packet on the slowest one if the difference between their RTTs is large enough. Even simple heuristics based on delay, transmission rate, and loss rate often perform better than the Lowest-RTT scheduler [196], [197].

In [198], the authors use a weighted round robin scheduler in combination with load balancing to overcome this difficulty; loss-aware scheduling is also a possibility [220]. Earliest Completion First (ECF) [199] is another scheduler that considers completion time as its main objective; it tries to reduce underutilization of flows by avoiding long idle periods, as they would cause CWND resets and consequent inefficiency in the capacity utilization. Decoupled Multipath Scheduler (DEMS) [200] uses the two paths to transmit the data outof-order: the data is sent from the first packet of a chunk on the first path, and from the last packet backwards, until the two flows meet and the chunk is fully downloaded.

In general, the optimal scheduling might not always send packets in-order, as sending future packets so that they arrive at the same time as the first (which is then sent later on a much faster flow) can be advantageous. The Slide Together Multipath Scheduler (STMS) [201] and Delay Aware Packet Scheduling (DAPS) [202] schedulers explicitly model this aspect, interleaving packets so that successive packets sent over different paths arrive at the same time. The FEC-based congestion control scheme in [215] also uses a similar strategy with good results.

The Blocking Estimation (BLEST) scheduler [203] is the first one in the literature to explicitly consider HoL; it tries to estimate which sub-flows are likely to cause it and dynamically adapts the scheduling to prevent it. Cross-layer scheduling is also a possibility; in [204], the authors schedule application-level objects on different paths to minimize webpage-loading times. The Q-aware scheduler [205] uses information from the network in a cross-layer fashion, directly considering

Protocol	Algorithm	Mechanism	Pros	Cons
MPTCP	LowRTT [196]	Lowest RTT first	Simple, widely deployed	Often inefficient
	Loss-aware [197]	Lowest RTT, weighted by loss rate	Works in wireless networks	Bad with asymmetric paths
	Weighted round robin [198]	Round robin, weighted by path performance	Optimal load balancing	Unknown fairness
	ECF [199]	Lowest completion time	Avoids idle flows	Can be hurt by HoL
	DEMS [200]	Forward on one path, backwards on the other	Object level delay	Requires data in blocks
	STMS [201], DAPS [202]	Delay modeling	Packets arrive in the correct order	RTT error sensitivity
	BLEST [203]	Explicit HoL minimization	Limits the HoL issue	RTT error sensitivity
	Application-aware [204]	Layer 7 delay minimization	Improved QoE	Requires application-level objects
	Q-aware [205]	Direct estimation of buffer occupancy	Limits the HoL issue	Requires network assistance
MPQUIC	Stream-aware [29]	Scheduling by stream	Improved webpage download times	Limited applicability outside HTTP/3
MP-DCCP	AOPS [195]	Delivery time and reliability modeling	Limits reordering	RTT error sensitivity
LEAP	DKF [30]	Deadline-based Kalman filter	Reliable latency with FEC	FEC

Table IV: Summary of the main presented scheduling algorithms

buffer occupancy. An experimental evaluation of several of the scheduling algorithms above is presented in [221].

We summarize the main features, advantages and drawbacks of most of the scheduling algorithms described above in Table IV.

3) Multipath TCP in data center networks: One of the original proposed use cases for MPTCP was the data center scenario [222]: since data centers often have complex topologies with multiple available paths between hosts, MPTCP is a natural solution. Large-scale simulations show that using MPTCP can improve load balancing, leading to fewer underutilized links and higher overall throughput, particularly in optical networks [223]. However, it is not immune to the Incast issue described in Sec. II-B [224]: whenever a client requests data from multiple servers at once, throughput collapses, even if MPTCP can actively relieve congested links. While the subflows from an MPTCP connection are not more aggressive than a single-path TCP flow, multiple MPTCP connections are not aware of each other. The solution proposed in [224] is to consider the number of existing flows when updating the congestion window of each subflow, combining an equally weighted coupling between flows to the subflow-level coupling described in Sec. V-A1.

The possibility of quickly retransmitting lost packets on less congested flows is another enhancement that has been proposed for the data center scenario [225]: this kind of technique is not needed for long-lived flows, but it can compensate for MPTCP's inability to steer short flows towards less congested paths and alleviate the Incast issue for this kind of frequent, short-lived traffic.

The benefits of MPTCP in data center networks can be increased when combined with SDN: network support can improve routing [226], leading to fewer congestion events. The performance benefit is even larger when MPTCP senders themselves are aware of the network situation and can dynamically add and remove subflows to avoid congestion [227].

#### B. Multipath in other protocols

The first alternative to TCP to have a multipath capability was SCTP [193]: the CMT extension increases throughput, but it has TCP friendliness issues. Several CMT-SCTP congestion control algorithms have been proposed [228], both using resource pooling concepts and following in MPTCP's footsteps; a full Markov model of congestion control is presented in [27]. An experimental comparison between MPTCP and CMT-SCTP was performed in [194]; the results showed that MPTCP has a slight performance advantage, but both protocols have issues with high-delay paths. For a more thorough analysis of CMT-SCTP congestion control algorithms, we refer the reader to the multipath congestion control survey we mentioned above [32], which examines it in detail.

More recent proposals add multipath capabilities to the QUIC protocol, paralleling the development of MPTCP: since QUIC is implemented in the user-space and not in the kernel, it is much easier to extend. Proposals for Multipath QUIC (MPQUIC) were recently advanced in [28], [229] and [29].

Both works develop similar aspects of the protocol, principally taking advantage of QUIC's features, like the 0-RTT connection establishment or its transparency to middleboxes. The great advantage of QUIC over TCP in the multipath domain is that it does not require in-order delivery on each sub-flow, sidestepping the HoL problem. For this reason, the design parameters and constraints for multipath congestion control in QUIC are not equivalent to those of MPTCP, and future work comparing them should take this into account. At the moment, research on MPQUIC is still very limited, as the multipath extension was only proposed in 2017. A first example of QUIC-specific work in multipath is given by [29], which proposes a stream-aware scheduler.

Finally, MP-DCCP is a recent draft [230] that extends DCCP to provide multipath capabilities. The peculiarities of DCCP are also present in its multipath version: data transfer is unreliable and out of order. Contrary to MPTCP, the communication establishment does not rely on a set initial path: as long as at least one functional path is available, the MP-DCCP connection can be established. The standard is still unfinished, and parts of the signaling and receiver-side reassembly procedures are yet to be defined.

Reordering is a major problem in MP-DCCP: even though single-path DCCP is also unreliable and does not guarantee inorder delivery, reordering is rare on single-path connections. In a multipath scenarios, reorderings can be much more common, and the use of a scheduler that prevents them as much as possible is recommended to prevent an excessive degradation of the application performance.

A possible multipath congestion control scheme is Flow Sharing Congestion Control (FSCC) [195], based on CCID

![](_page_17_Figure_2.jpeg)

Figure 14: Performance of well-known multipath protocols in the throughput/reliability plane in a multipath scenario [Wi-Fi+LTE]. Reprinted, with permission, from [30].

2. It adds a "helping state" to each flow's congestion control state, which is triggered when another path is congested. After the congested path's CWND is reduced, the helping path's is increased accordingly to increase the overall throughput stability. The authors of FSCC also propose a predictive scheduler called Adaptive Order Prediction Scheduling (AOPS) as a way to deal with the reordering issue: the scheme considers the predicted delivery time of packets, as well as the reliability of each path, to optimize the schedule for in-order delivery. Reliability can also be adapted to the type of data being sent, protecting more important data so that it is delivered on the most reliable path, even if this increases the latency.

There are also some natively multipath protocols: LEAP [30] is a protocol developed specifically to exploit the multipath scenario in order to guarantee low latency to applications. It uses a capacity-based congestion control mechanism on each sub-flow, and adaptively sends FEC to compensate for capacity estimation errors, exploiting path diversity to guarantee that packets will reach the receiver with a strictly bounded latency. Since it runs in user-space over UDP, it does not require retransmission on each sub-flow: lost or late packets that are recovered using the FEC protection on the other paths do not need to be retransmitted at all.

Fig. 14 shows the relative performance of several congestion control mechanisms in a wireless scenario with two paths, one over LTE and one over an office Wi-Fi: while MPTCP BALIA performs far better than the uncoupled CUBIC because it handles the HoL problem far better, it does not provide any latency guarantees. Less aggressive protocols such as Verus, Vegas, and Sprout do better in terms of latency, but they still do not manage to go below a lateness rate of 10%. LEAP can control the trade-off between the strictness of the latency constraint and the throughput by dynamically adjusting the coding rate, and it is the only protocol that manages to violate the latency requirement less than 1% of the time.

#### C. Open Challenges and Research Directions

As discussed in Sec. III and in the previous paragraphs, the exploitation of multiple paths at the transport layer is a promising research trend, made possible by the advanced capabilities of recent communication devices, such as smartphones capable of connecting to the network over multiple wireless technologies (usually cellular and Wi-Fi). The field is still in its infancy, and there are unsolved issues such as fairness and HoL blocking: some of these issues are caused by the backward compatibility requirement and can be mitigated by using flexible protocols such as MPQUIC, but other problems are more fundamental.

The close coupling of scheduling and congestion control, which have significant interactions with hard-to-model effects, is one of the biggest issues for multipath protocols, and it affects all protocols: all the issues described in Sec. II are exacerbated by the additional layer of adaptation. Designing an efficient integrated multipath scheduling and congestion control scheme is very challenging, and will probably require the use of FEC and extensive cross-layer support. Finally, the use of multipath protocols in wireless and data center networks requires specific adaptations, and there are several research challenges that still have to be solved.

## VI. CONCLUSIONS

TCP has been the *de facto* standard transport protocol for years, but, despite its wide adoption, it presents suboptimal performance in a number of use cases and scenarios. Moreover, new emerging technologies (e.g., mmWave communications) and requirements (e.g., those for ultra-low latency VR streaming) are overtaking the performance that TCP can achieve. Therefore, the research related to the transport layer has seen a renewed interest in the last few years. In this paper, we reviewed the main results related to these efforts. In particular, we analyzed three main research areas: transport protocols in general, congestion control, and multipath transport.

First, we discussed the main features of three proposed protocols, namely QUIC, SCTP and DCCP, which represent significant evolutions or alternatives to the widely used TCP. These protocols have several improvements related to the acknowledgment mechanism, the connection establishment, the management of the connection life cycle, and the embedding of the cryptographic stack in the protocol itself. SCTP and DCCP have been standardized by the IETF, but have not reached a wide adoption in the Internet. QUIC, instead, has not been fully standardized yet, but is implemented in the user space, thus can be deployed without the need to upgrade the operating system. Besides, some early results show its promising performance, at the cost of a higher computational load with respect to kernel-based solutions. QUIC can be seen as a possible enabler of a flexible deployment of congestion control to target different scenarios and use cases.

Then, we reviewed a selection of recently proposed congestion control algorithms for TCP and, in general, for congestion-aware transport protocols. We analyzed different approaches, from evolutions of traditional loss-based mechanisms to new delay-based or hybrid proposals and machine learning strategies.

Finally, another promising trend is represented by the usage of multiple paths at the transport layer, to provide macro diversity and, possibly, increase the throughput and reliability. A first solution is MPTCP, an extension of TCP that dispatches packets over multiple subflows when multiple network interfaces are available. However, while maintaining TCP compatibility with respect to middleboxes by design, it suffers from HoL issues (depending on the scheduler implementation) and fairness with respect to single path TCP. Multi-path extensions of other protocols (e.g., SCTP, QUIC, and DCCP) have also been proposed but are still open research areas.

The research related to several of the topics presented in this survey is still ongoing, and, given the fast update rate of communication technologies, the interest in innovation at the transport layer will hardly fade away. A promising research direction is related to the coupling of network slicing and TCP fairness: slices can be instantiated by network operators to protect delay-based flows against less fair algorithms. Similarly, adaptive congestion control algorithms are still being studied, with the aim of targeting the optimal operating point that minimizes the latency while maximizing the data rate. Finally, in the multipath domain, standardization efforts are still ongoing, and new efficient and fair multi-link scheduling and congestion control algorithms will need to be identified.

#### REFERENCES

- J. Postel, "Transmission Control Protocol," IETF, RFC 793, Sep. 1981. [Online]. Available: https://rfc-editor.org/rfc/rfc793.txt
- [2] —, "Internet Protocol," IETF, RFC 791, Sep. 1981. [Online]. Available: https://rfc-editor.org/rfc/rfc791.txt
- [3] —, "User Datagram Protocol," RFC 768, Aug. 1980. [Online]. Available: https://rfc-editor.org/rfc/rfc768.txt
- [4] Cisco, "Cisco Visual Networking Index, forecast and methodology, 2016–2021," Tech. Rep., June 2017.
- [5] 3GPP, "NR and NG-RAN Overall Description Rel. 15," TS 38.300, Jan. 2018.
- [6] S. Parkvall, E. Dahlman, A. Furuskar, and M. Frenne, "NR: The new 5G Radio Access Technology," *IEEE Communications Standards Magazine*, vol. 1, no. 4, pp. 24–30, Dec 2017.
- [7] B. Bellalta, "IEEE 802.11ax: High-efficiency WLANs," *IEEE Wireless Communications*, vol. 23, no. 1, pp. 38–46, February 2016.
- [8] Y. Ghasempour, C. R. C. M. da Silva, C. Cordeiro, and E. W. Knightly, "IEEE 802.11ay: Next-generation 60 GHz communication for 100 Gb/s Wi-Fi," *IEEE Communications Magazine*, vol. 55, no. 12, pp. 186–192, Dec. 2017.
- [9] Q. D. Coninck, M. Baerts, B. Hesmans, and O. Bonaventure, "Observing real smartphone applications over multipath TCP," *IEEE Communications Magazine*, vol. 54, no. 3, pp. 88–93, Mar. 2016.
- [10] R. Alvizu, G. Maier, N. Kukreja, A. Pattavina, R. Morro, A. Capello, and C. Cavazzoni, "Comprehensive survey on T-SDN: Software-Defined Networking for transport networks," *IEEE Communications Surveys and Tutorials*, vol. 19, no. 4, pp. 2232–2283, Jan. 2017.
- [11] X. Liu, A. Sridharan, S. Machiraju, M. Seshadri, and H. Zang, "Experiences in a 3G network: Interplay between the wireless channel and applications," in *14th ACM International Conference on Mobile Computing and Networking (MobiCom)*, San Francisco, California, USA, Sep. 2008, pp. 211–222.
- [12] R. Stewart, "Stream Control Transport Protocol," IETF, RFC 6582, Sep. 2007. [Online]. Available: https://rfc-editor.org/rfc/rfc6582.txt
- [13] P. Srisuresh, J. Kuthan, J. Rosenberg, A. Molitor, and A. Rayhan, "Middlebox communication architecture and framework," IETF, RFC 3303, Aug. 2002. [Online]. Available: https://rfc-editor.org/rfc/rfc3303. txt

- [14] B. Carpenter and S. Brim, "Middleboxes: Taxonomy and issues," IETF, RFC 3234, Feb. 2002. [Online]. Available: https://rfc-editor.org/ rfc/rfc3234.txt
- [15] K. Edeline and B. Donnet, "A first look at the prevalence and persistence of middleboxes in the wild," in 29th IEEE International Teletraffic Congress (ITC), vol. 1, Genoa, Italy, Sep. 2017, pp. 161– 168.
- [16] G. Papastergiou, G. Fairhurst, D. Ros, A. Brunstrom, K. Grinnemo, P. Hurtig, N. Khademi, M. Tüxen, M. Welzl, D. Damjanovic, and S. Mangiante, "De-ossifying the Internet transport layer: A survey and future perspectives," *IEEE Communications Surveys and Tutorials*, vol. 19, no. 1, pp. 619–639, Feb. 2017.
- [17] H.-S. Park, J.-Y. Lee, and B.-C. Kim, "TCP performance issues in LTE networks," in *International Conference on ICT Convergence (ICTC)*. Seoul, South Korea: IEEE, Sep. 2011, pp. 493–496.
- [18] M. Zhang, M. Polese, M. Mezzavilla, J. Zhu, S. Rangan, S. Panwar, and a. M. Zorzi, "Will TCP work in mmWave 5G cellular networks?" *IEEE Communications Magazine*, vol. 57, no. 1, pp. 65–71, Jan. 2019.
- [19] M. Scharf and S. Kiesel, "Head-of-line blocking in TCP and SCTP: Analysis and measurements," in *IEEE Global Communications Conference (GLOBECOM)*, San Francisco, California, USA, Nov. 2006, pp. 1–5.
- [20] J. Gettys and K. Nichols, "Bufferbloat: Dark buffers in the Internet," ACM Queue, vol. 9, no. 11, pp. 40:40–40:54, Nov. 2011.
- [21] N. Cardwell, Y. Cheng, C. S. Gunn, S. H. Yeganeh, and V. Jacobson, "BBR: Congestion-based congestion control," ACM Queue, vol. 14, no. 5, pp. 20–53, Sep. 2016.
- [22] M. Hock, F. Neumeister, M. Zitterbart, and R. Bless, "TCP LoLa: Congestion control for low latencies and high throughput," in 42nd IEEE Conference on Local Computer Networks (LCN), Singapore, Oct. 2017, pp. 215–218.
- [23] E. Kohler, M. Handley, and S. Floyd, "Datagram Congestion Control Protocol (DCCP)," IETF, RFC 4340, March 2006. [Online]. Available: http://www.rfc-editor.org/rfc/rfc4340.txt
- [24] J. Iyengar and M. Thomson, "QUIC: A UDP-based multiplexed and secure transport," IETF, Working Draft: draft-ietf-quic-transport-08, Dec. 2017. [Online]. Available: https://tools.ietf.org/id/draft-ietf-quic-transport-08.txt
- [25] A. Ford, C. Raiciu, M. Handley, S. Barre, and J. Iyengar, "Architectural guidelines for multipath TCP development," IETF, RFC 6182, Mar. 2011. [Online]. Available: https://rfc-editor.org/rfc/rfc6182.txt
- [26] C. Raiciu, M. J. Handley, and D. Wischik, "Coupled congestion control for multipath transport protocols," IETF, RFC 6356, Oct. 2011. [Online]. Available: https://rfc-editor.org/rfc/rfc6356.txt
- [27] T. D. Wallace and A. Shami, "Concurrent multipath transfer using SCTP: Modelling and congestion window management," *IEEE Transactions on Mobile Computing*, vol. 13, no. 11, pp. 2510–2523, Nov. 2014.
- [28] Q. D. Coninck and O. Bonaventure, "Multipath extension for QUIC," IETF, Working Draft: deconinck-multipath-quic-00, Oct. 2017. [Online]. Available: https://datatracker.ietf.org/doc/html/ draft-deconinck-multipath-quic-00
- [29] T. Viernickel, A. Froemmgen, A. Rizk, B. Koldehofe, and R. Steinmetz, "Multipath QUIC: A deployable multipath transport protocol," in *IEEE International Conference on Communications (ICC)*, May 2018, pp. 1–7.
- [30] F. Chiariotti, S. Kucera, A. Zanella, and H. Claussen, "Analysis and design of a latency control protocol for multi-path data delivery with pre-defined QoS guarantees," *IEEE/ACM Transactions on Networking*, vol. 27, no. 3, pp. 1165–1178, Jun. 2019.
- [31] M. Li, A. Lukyanenko, Z. Ou, A. Ylä-Jääski, S. Tarkoma, M. Coudron, and S. Secci, "Multipath transmission for the internet: A survey," *IEEE Communications Surveys and Tutorials*, vol. 18, no. 4, pp. 2887–2925, Jun. 2016.
- [32] C. Xu, J. Zhao, and G. M. Muntean, "Congestion control design for multipath transport protocols: A survey," *IEEE Communications Surveys and Tutorials*, vol. 18, no. 4, pp. 2948–2969, Apr. 2016.
- [33] R. Rojas-Cessa, Y. Kaymak, and Z. Dong, "Schemes for fast transmission of flows in data center networks," *IEEE Communications Surveys* and Tutorials, vol. 17, no. 3, pp. 1391–1422, Aug. 2015.
- [34] T. Lukaseder, L. Bradatsch, B. Erb, R. W. Van Der Heijden, and F. Kargl, "A comparison of TCP congestion control algorithms in 10G networks," in *41st IEEE Conference on Local Computer Networks* (*LCN*), Dubai, UAE, Nov. 2016, pp. 706–714.
- [35] Y. Gong, D. Rossi, C. Testa, S. Valenti, and M. D. Täht, "Fighting the bufferbloat: On the coexistence of AQM and low priority congestion

control," in *IEEE Conference on Computer Communications (INFO-COM)*, Turin, Italy, Apr. 2013, pp. 3291–3296.

- [36] Z. Liu, J. Sun, S. Hu, and X. Hu, "An adaptive AQM algorithm based on a novel information compression model," *IEEE Access*, vol. 6, pp. 31 180–31 190, Jun. 2018.
- [37] C. A. Grazia, N. Patriciello, M. Klapez, and M. Casoni, "A crosscomparison between TCP and AQM algorithms: Which is the best couple for congestion control?" in *14th IEEE International Conference* on Telecommunications (ConTEL), Zagreb, Croatia, Jun. 2017, pp. 75– 82.
- [38] K. Nichols and V. Jacobson, "Controlling queue delay," Communications of the ACM, vol. 55, no. 7, pp. 42–50, Jul. 2012.
- [39] R. Amrutha and V. Nithya, "Curbing of TCP Incast in data center networks," in 4th IEEE International Conference on Reliability, Infocom Technologies and Optimization (ICRITO) (Trends and Future Directions), Noida, India, Sep. 2015, pp. 1–5.
- [40] J. Dean and S. Ghemawat, "MapReduce: Simplified data processing on large clusters," in Usenix 6th Symposium on Operating System Design and Implementation (OSDI), San Francisco, California, USA, Dec. 2004, pp. 137–150.
- [41] L. Page, S. Brin, R. Motwani, and T. Winograd, "The PageRank citation ranking: Bringing order to the web," Stanford InfoLab, Tech. Rep., Nov. 1999.
- [42] W. Chen, F. Ren, J. Xie, C. Lin, K. Yin, and F. Baker, "Comprehensive understanding of TCP Incast problem," in *IEEE Conference on Computer Communications (INFOCOM)*, Hong Kong, China, Apr. 2015, pp. 1688–1696.
- [43] Y. Ren, Y. Zhao, P. Liu, K. Dou, and J. Li, "A survey on TCP Incast in data center networks," *International Journal of Communication Systems*, vol. 27, no. 8, pp. 1160–1172, Aug. 2014.
- [44] J. Zhang, F. Ren, L. Tang, and C. Lin, "Modeling and solving TCP Incast problem in data center networks," *IEEE Transactions on Parallel* and Distributed Systems, vol. 26, no. 2, pp. 478–491, Feb. 2015.
- [45] M. Alizadeh, A. Greenberg, D. A. Maltz, J. Padhye, P. Patel, B. Prabhakar, S. Sengupta, and M. Sridharan, "Data Center TCP (DCTCP)," in ACM Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications (SIGCOMM). New Delhi, India: ACM, 2010, pp. 63–74.
- [46] J. Hwang, J. Yoo, and N. Choi, "IA-TCP: A rate based Incast-avoidance algorithm for TCP in data center networks," in *IEEE International Conference on Communications (ICC)*, Ottawa, Canada, Jun. 2012, pp. 1292–1296.
- [47] J. Zhang, F. Ren, and C. Lin, "Survey on transport control in data center networks," *IEEE Network*, vol. 27, no. 4, pp. 22–26, Jul. 2013.
- [48] J. Iyengar and I. Swett, "QUIC loss detection and congestion control," IETF, Working Draft: draft-ietf-quic-recovery-08, Dec. 2017. [Online]. Available: https://tools.ietf.org/id/draft-ietf-quic-recovery-08.txt
- [49] M. Seemann and L. Clemente. (2019, Apr.) QUIC implementation in Go. [Online]. Available: github.com/lucas-clemente/quic-go
- [50] D. Tikhonov and B. Prodoehl. (2017, Sep.) LiteSpeed QUIC. [Online]. Available: github.com/litespeedtech/lsquic-client
- [51] J. Lee, B. Hong *et al.* (2016, Sep.) libquic. [Online]. Available: github.com/devsisters/libquic
- [52] R. Shade *et al.* (2016, Oct.) proto-quic. [Online]. Available: github.com/google/proto-quic
- [53] G. Fairhurst, "The Datagram Congestion Control Protocol (DCCP) Service Codes," IETF, RFC 5595, Sep. 2009. [Online]. Available: https://rfc-editor.org/rfc/rfc5595.txt
- [54] —, "Datagram Congestion Control Protocol (DCCP) simultaneousopen technique to facilitate NAT/middlebox traversal," IETF, RFC 5596, Sep. 2009. [Online]. Available: https://rfc-editor.org/rfc/rfc5596. txt
- [55] S. Floyd, J. Mahdavi, M. Mathis, and D. A. Romanow, "TCP Selective Acknowledgment options," IETF, RFC 2018, Oct. 1996. [Online]. Available: https://rfc-editor.org/rfc/rfc2018.txt
- [56] M. Belshe, M. Thomson, and R. Peon, "Hypertext Transfer Protocol version 2," IETF, RFC 7540, May 2015. [Online]. Available: https://rfc-editor.org/rfc/rfc7540.txt
- [57] S. McQuistin, C. Perkins, and M. Fayed, "TCP Hollywood: An unordered, time-lined, TCP for networked multimedia applications," in *IEEE/IFIP Networking Conference (IFIP Networking)*, Vienna, Austria, May 2016, pp. 422–430.
- [58] Y. Cui, T. Li, C. Liu, X. Wang, and M. Kühlewind, "Innovating transport with QUIC: Design approaches and research challenges," *IEEE Internet Computing*, vol. 21, no. 2, pp. 72–76, Mar. 2017.

- [59] H. Balakrishnan, V. N. Padmanabhan, S. Seshan, and R. H. Katz, "A comparison of mechanisms for improving TCP performance over wireless links," *IEEE/ACM Transactions on Networking*, vol. 5, no. 6, pp. 756–769, Dec 1997.
- [60] G. Xylomenos, G. C. Polyzos, P. Mahonen, and M. Saaranen, "TCP performance issues over wireless links," *IEEE Communications Magazine*, vol. 39, no. 4, pp. 52–58, April 2001.
- [61] J. Griner, J. L. Border, M. Kojo, Z. D. Shelby, and G. Montenegro, "Performance enhancing proxies intended to mitigate link-related degradations," IETF, RFC 3135, Jun. 2001. [Online]. Available: https://rfc-editor.org/rfc/rfc3135.txt
- [62] K. Liu and J. Y. Lee, "On improving TCP performance over mobile data networks," *IEEE Transactions on Mobile Computing*, vol. 15, no. 10, pp. 2522–2536, Oct. 2016.
- [63] M. Maity, B. Raman, and M. Vutukuru, "TCP download performance in dense WiFi scenarios: Analysis and solution," *IEEE Transactions on Mobile Computing*, vol. 16, no. 1, pp. 213–227, Jan. 2017.
- [64] B. Sardar and D. Saha, "A survey of TCP enhancements for lasthop wireless networks," *IEEE Communications Surveys and Tutorials*, vol. 8, no. 3, pp. 20–34, Jul. 2006.
- [65] S. Rangan, T. S. Rappaport, and E. Erkip, "Millimeter-wave cellular wireless networks: Potentials and challenges," *Proceedings of the IEEE*, vol. 102, no. 3, pp. 366–385, Mar. 2014.
- [66] M. Zhang, M. Mezzavilla, R. Ford, S. Rangan, S. S. Panwar, E. Mellios, D. Kong, A. R. Nix, and M. Zorzi, "Transport Layer Performance in 5G mmWave Cellular," in *IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*, San Francisco, California, USA, Apr. 2016, pp. 730–735.
- [67] M. Polese, R. Jana, and M. Zorzi, "TCP and MP-TCP in 5G mmWave networks," *IEEE Internet Computing*, vol. 21, no. 5, pp. 12–19, Sep. 2017.
- [68] H. Jiang, Z. Liu, Y. Wang, K. Lee, and I. Rhee, "Understanding bufferbloat in cellular networks," in ACM SIGCOMM Workshop on Cellular Networks: Operations, Challenges, and Future Design (Cell-Net), Helsinki, Finland, Aug. 2012, pp. 1–6.
- [69] A. M. Kakhki, S. Jero, D. Choffnes, C. Nita-Rotaru, and A. Mislove, "Taking a long look at QUIC: An approach for rigorous evaluation of rapidly evolving transport protocols," in ACM Internet Measurement Conference (IMC), London, United Kingdom, Nov. 2017, pp. 290–303.
- [70] A. Langley, A. Riddoch, A. Wilk, A. Vicente, C. Krasić, D. Zhang, F. Yang, F. Kouranov, I. Swett, J. Iyengar, J. Bailey, J. Dorfman, J. Roskind, J. Kulik, P. Westin, R. Tenneti, R. Shade, R. Hamilton, V. Vasiliev, W.-T. Chang, and Z. Shi, "The QUIC transport protocol: Design and Internet-scale deployment," in ACM Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications (SIGCOMM). Los Angeles, California, USA: ACM, Aug. 2017, pp. 183–196.
- [71] E. Rescorla, "The Transport Layer Security (TLS) protocol version 1.3," IETF, RFC 8446, Aug. 2018. [Online]. Available: https: //rfc-editor.org/rfc/rfc8446.txt
- [72] G. Carlucci, L. D. Cicco, and S. Mascolo, "HTTP over UDP: An experimental investigation of QUIC," in *30th Annual ACM Symposium* on Applied Computing (SAC), Salamanca, Spain, Apr. 2015, pp. 609– 614.
- [73] P. Megyesi, Z. Krämer, and S. Molnár, "How quick is QUIC?" in *IEEE International Conference on Communications (ICC)*, Kuala Lumpur, Malaysia, May 2016, pp. 1–6.
- [74] J. Rüth, I. Poese, C. Dietzel, and O. Hohlfeld, "A first look at QUIC in the wild," in *Passive and Active Measurement Conference (PAM)*, Berlin, Germany, 2018, pp. 255–268.
- [75] A. D. Biasio, F. Chiariotti, M. Polese, A. Zanella, and M. Zorzi, "A QUIC implementation for ns-3," in *Workshop on ns-3 (WNS3)*, Florence, Italy, Jun. 2019, pp. 1–7.
- [76] P. Natarajan, F. Baker, P. D. Amer, and J. T. Leighton, "SCTP: What, why, and how," *IEEE Internet Computing*, vol. 13, no. 5, pp. 81–85, Sep. 2009.
- [77] R. Stewart and C. Metz, "SCTP: New transport protocol for TCP/IP," *IEEE Internet Computing*, vol. 5, no. 6, pp. 64–69, Nov. 2001.
- [78] S. Fu and M. Atiquzzaman, "SCTP: state of the art in research, products, and technical challenges," *IEEE Communications Magazine*, vol. 42, no. 4, pp. 64–76, April 2004.
- [79] A. L. Caro, J. R. Iyengar, P. D. Amer, S. Ladha, G. J. Heinz, and K. C. Shah, "SCTP: A proposed standard for robust Internet data transport," *IEEE Computer*, vol. 36, no. 11, pp. 56–63, Nov. 2003.
- [80] P. T. Conrad, G. J. Heinz, A. L. Caro, P. D. Amer, and J. Fiore, "SCTP in battlefield networks," in *Military Communications Conference (MIL-COM)*, vol. 1, McLean, Virginia, USA, oct 2001, pp. 289–295.

21

- [81] A. Jungmaier, M. Schopp, and M. Tuxen, "Performance evaluation of the simple control transmission protocol (SCTP)," in *IEEE Conference* on High Performance Switching and Routing (HPSR), Heidelberg, Germany, Jun. 2000, pp. 141–148.
- [82] J. Shi, Y. Jin, H. Huang, and D. Zhang, "Experimental performance studies of SCTP in wireless access networks," in *International Conference on Communication Technology (ICCT)*, vol. 1, Beijing, China, Apr. 2003, pp. 392–395.
- [83] G. Camarillo, R. Kantola, and H. Schulzrinne, "Evaluation of transport protocols for the Session Initiation Protocol," *IEEE Network*, vol. 17, no. 5, pp. 40–46, Sep. 2003.
- [84] E. Kohler, M. Handley, and S. Floyd, "Designing dccp: Congestion control without reliability," in ACM Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications (SIGCOMM), Pisa, Italy, 2006, pp. 27–38.
- [85] L. Eggert, G. Fairhurst, and G. Shepherd, "UDP usage guidelines," IETF, RFC 8085, Mar. 2017. [Online]. Available: https://rfc-editor. org/rfc/rfc8085.txt
- [86] S. Floyd, M. Handley, and E. Kohler, "Problem Statement for the Datagram Congestion Control Protocol (DCCP)," IETF, RFC 4336, Mar. 2006. [Online]. Available: http://www.rfc-editor.org/rfc/rfc4336. txt
- [87] K. Ramakrishnan, S. Floyd, and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP," IETF, RFC 3168, Sep. 2001. [Online]. Available: http://www.rfc-editor.org/rfc/rfc3168.txt
- [88] A. Gurtov, T. Henderson, S. Floyd, and Y. Nishida, "The NewReno modification to TCP's Fast Recovery algorithm," IETF, RFC 6582, Apr. 2012. [Online]. Available: https://rfc-editor.org/rfc/rfc6582.txt
- [89] L. Xu, K. Harfoush, and I. Rhee, "Binary increase congestion control (BIC) for fast long-distance networks," in *IEEE Conference on Computer Communications (INFOCOM)*, vol. 4, Hong Kong, China, Mar. 2004, pp. 2514–2524.
- [90] S. Ha, I. Rhee, and L. Xu, "CUBIC: A new TCP-friendly high-speed TCP variant," ACM Operating Systems Review, vol. 42, no. 5, pp. 64– 74, Jul. 2008.
- [91] A. Abdelsalam, M. Luglio, C. Roseti, and F. Zampognaro, "TCP Wave: a new reliable transport approach for future Internet," *Computer Networks*, vol. 112, pp. 122–143, Jan. 2017.
- [92] L. S. Brakmo, S. W. O'Malley, and L. L. Peterson, "TCP Vegas: New techniques for congestion detection and avoidance," ACM Computer Communication Review, vol. 24, no. 4, pp. 24–35, Oct. 1994.
- [93] J. Sing and B. Soh, "TCP New Vegas: Performance evaluation and validation," in 11th IEEE Symposium on Computers and Communications (ISCC), Cagliari, Italy, Jun. 2006, pp. 541–546.
- [94] Y. Zaki, T. Pötsch, J. Chen, L. Subramanian, and C. Görg, "Adaptive congestion control for unpredictable cellular networks," in ACM Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications (SIGCOMM), London, UK, 2015, pp. 509–522.
- [95] P. Goyal, A. Narayan, F. Cangialosi, D. Raghavan, S. Narayana, M. Alizadeh, and H. Balakrishnan, "Elasticity Detection: A Building Block for Delay-Sensitive Congestion Control," *arXiv* preprint arXiv:1802.08730, pp. 1–17, Jul. 2018. [Online]. Available: https://arxiv.org/abs/1802.08730
- [96] S. Shalunov, G. Hazel, J. Iyengar, and M. Kuehlewind, "Low Extra Delay Background Transport (LEDBAT)," IETF, RFC 6817, Dec. 2012. [Online]. Available: https://rfc-editor.org/rfc/rfc6817.txt
- [97] S. Mascolo, C. Casetti, M. Gerla, M. Y. Sanadidi, and R. Wang, "TCP Westwood: Bandwidth estimation for enhanced transport over wireless links," in 7th ACM conference on Mobile Computing and Networking (MobiCom), Rome, Italy, Jul. 2001, pp. 287–297.
- [98] K. Winstein, A. Sivaraman, and H. Balakrishnan, "Stochastic forecasts achieve high throughput and low delay over cellular networks," in 10th USENIX Conference on Networked Systems Design and Implementation (NSDI), Lombard, Illinois, USA, Apr. 2013, pp. 459–472.
- [99] K. Tan, J. Song, Q. Zhang, and M. Sridharan, "A Compound TCP approach for high-speed and long distance networks," in 25th IEEE International Conference on Computer Communications (INFOCOM), Barcelona, Spain, Apr. 2006, pp. 1–12.
- [100] S. Liu, T. Başar, and R. Srikant, "TCP-Illinois: A loss-and delay-based congestion control algorithm for high-speed networks," *Performance Evaluation*, vol. 65, no. 6–7, pp. 417–440, Jun. 2008.
- [101] C. P. Fu and S. C. Liew, "TCP Veno: TCP enhancement for transmission over wireless access networks," *IEEE Journal on selected areas in communications*, vol. 21, no. 2, pp. 216–228, Feb 2003.

- [102] K. Winstein and H. Balakrishnan, "TCP Ex Machina: Computergenerated congestion control," ACM Computer Communication Review, vol. 43, no. 4, pp. 123–134, Aug. 2013.
- [103] A. Sivaraman, K. Winstein, P. Thaker, and H. Balakrishnan, "An experimental study of the learnability of congestion control," ACM Computer Communication Review, vol. 44, no. 4, pp. 479–490, Aug. 2014.
- [104] M. Dong, Q. Li, D. Zarchy, P. B. Godfrey, and M. Schapira, "PCC: Re-architecting congestion control for consistent high performance," in *12th USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, Oakland, Califonia, USA, May 2015, pp. 395–408.
- [105] X. Nie, Y. Zhao, Z. Li, G. Chen, K. Sui, J. Zhang, Z. Ye, and D. Pei, "Dynamic TCP Initial Windows and Congestion Control Schemes Through Reinforcement Learning," *IEEE Journal on Selected Areas* in Communications, vol. 37, no. 6, pp. 1231–1247, June 2019.
- [106] W. Li, F. Zhou, K. R. Chowdhury, and W. M. Meleis, "QTCP: Adaptive Congestion Control with Reinforcement Learning," *IEEE Transactions* on Network Science and Engineering, pp. 1–13, May 2018, Early Access.
- [107] V.-H. Tran and O. Bonaventure, "Beyond socket options: making the Linux TCP stack truly extensible," *arXiv preprint arXiv:1901.01863*, pp. 1–9, Jan. 2019. [Online]. Available: https://arxiv.org/pdf/1901. 01863.pdf
- [108] R. Ando, T. Murase, and M. Oguchi, "Characteristics of QoSguaranteed TCP on real mobile terminal in wireless LAN," in *IEEE International Communications Quality and Reliability Workshop (CQR)*. Naples, Florida, USA: IEEE, May 2011, pp. 1–6.
- [109] P. M. Mohan, D. M. Divakaran, and M. Gurusamy, "Performance study of TCP flows with QoS-supported OpenFlow in data center networks," in 19th IEEE International Conference on Networks (ICON). Singapore: IEEE, Apr. 2013, pp. 1–6.
- [110] J. Nagle, "Congestion control in IP/TCP internetworks," IETF, RFC 896, Jan. 1984. [Online]. Available: https://rfc-editor.org/rfc/rfc896.txt
- [111] L. A. Grieco and S. Mascolo, "Performance evaluation and comparison of Westwood+, New Reno, and Vegas TCP congestion control," ACM Computer Communication Review, vol. 34, no. 2, pp. 25–38, Apr. 2004.
- [112] J. Sing and B. Soh, "TCP New Vegas: Improving the performance of TCP Vegas over high latency links," in *4th IEEE International Symposium on Network Computing and Application (NCA)*, Cambridge, Massachusetts, USA, Jul. 2005, pp. 73–82.
- [113] S. ullah Lar and X. Liao, "An initiative for a classified bibliography on TCP/IP congestion control," *Journal of Network and Computer Applications*, vol. 36, no. 1, pp. 126–133, Jan. 2013.
- [114] M. Allman, V. Paxson, and E. Blanton, "TCP congestion control," IETF, RFC 5681, Sep. 2009. [Online]. Available: https://rfc-editor.org/ rfc/rfc5681.txt
- [115] E. Blanton, M. Allman, L. Wang, I. Jarvinen, M. Kojo, and Y. Nishida, "A conservative loss recovery algorithm based on Selective Acknowledgment (SACK) for TCP," IETF, RFC 6675, Aug. 2012. [Online]. Available: https://rfc-editor.org/rfc/rfc6675.txt
- [116] D. Borman, R. T. Braden, V. Jacobson, and R. Scheffenegger, "TCP extensions for high performance," IETF, RFC 7323, Sep. 2014. [Online]. Available: https://rfc-editor.org/rfc/rfc7323.txt
- [117] M. Šošić and V. Stojanović, "Resolving poor TCP performance on high-speed long distance links – Overview and comparison of BIC, CUBIC and Hybla," in 11th IEEE International Symposium on Intelligent Systems and Informatics (SISY), Subotica, Serbia, Sep. 2013, pp. 325–330.
- [118] M. Ahmad, A. B. Ngadi, A. Nawaz, U. Ahmad, T. Mustafa, and A. Raza, "A survey on TCP CUBIC variant regarding performance," in 15th International Multitopic Conference (INMIC), Islamabad, Pakistan, Dec. 2012, pp. 409–412.
- [119] I. Rhee, L. Xu, S. Ha, A. Zimmermann, L. Eggert, and R. Scheffenegger, "CUBIC for fast long-distance networks," IETF, draft-ietg-tcpm-cubic-07, Nov. 2017. [Online]. Available: https: //tools.ietf.org/id/draft-ietg-tcpm-cubic-07
- [120] S. Floyd and E. Kohler, "Profile for Datagram Congestion Control Protocol (DCCP) Congestion Control ID 2: TCP-like congestion control," IETF, RFC 4341, Mar. 2006. [Online]. Available: http://www.rfc-editor.org/rfc/rfc4341.txt
- [121] S. Floyd, E. Kohler, and J. Padhye, "TCP-Friendly Rate Control (TFRC): Protocol specification," IETF, RFC 5348, Sep. 2008. [Online]. Available: http://www.rfc-editor.org/rfc/rfc5348.txt
- [122] S. Floyd and E. Kohler, "Profile for Datagram Congestion Control Protocol (DCCP) Congestion ID 4: TCP-Friendly Rate Control for Small Packets (TFRC-SP)," IETF, RFC 5622, Aug. 2009. [Online]. Available: http://www.rfc-editor.org/rfc/rfc5622.txt

- [123] M. Schier and M. Welzl, "Using DCCP: Issues and improvements," in 20th IEEE International Conference on Network Protocols (ICNP). Austin, Texas, USA: IEEE, Oct. 2012, pp. 1–9.
- [124] A. Kumar, L. Jacob, and A. L. Ananda, "SCTP vs TCP: Performance comparison in MANETs," in 29th IEEE International Conference on Local Computer Networks (LCN), Tampa, Florida, USA, Nov. 2004, pp. 431–432.
- [125] I. Ahmed, O. Yasuo, and K. Masanori, "Improving performance of SCTP over broadband high latency networks," in 28th IEEE International Conference on Local Computer Networks (LCN), Bonn, Germany, Oct 2003, pp. 644–645.
- [126] C. Roseti and E. Kristiansen, "TCP Noordwijk: TCP-based transport optimized for web traffic in satellite networks," in 26th AIAA International Communications Satellite Systems Conference (ICSSC), San Diego, California, USA, Jun. 2008, pp. 1–12.
- [127] L. S. Brakmo and L. L. Peterson, "TCP Vegas: end to end congestion avoidance on a global Internet," *IEEE Journal on Selected Areas in Communications*, vol. 13, no. 8, pp. 1465–1480, Oct. 1995.
- [128] K. N. Srijith, L. Jacob, and A. L. Ananda, "TCP Vegas-A: solving the fairness and rerouting issues of TCP Vegas," in *IEEE International Performance, Computing, and Communications Conference (IPCC)*, Phoenix, Arizona, USA, Apr. 2003, pp. 309–316.
- [129] W. Zhou, W. Xing, Y. Wang, and J. Zhang, "TCP Vegas-V: Improving the performance of TCP Vegas," in *IEEE International Conference on Automatic Control and Artificial Intelligence (ACAI)*, Xiamen, China, Mar. 2012, pp. 2034–2039.
- [130] Y. Guo, X. Yang, R. Wang, and J. Sun, "TCP Adaptive Vegas: Improving of TCP Vegas algorithm," in *IEEE International Conference on Information Science, Electronics and Electrical Engineering (ISEEE)*, vol. 1, Sapporo, Japan, Apr. 2014, pp. 126–130.
- [131] L. Ding, X. Wang, Y. Xu, W. Zhang, and W. Chen, "Vegas-W: An Enhanced TCP-Vegas for Wireless Ad Hoc Networks," in *IEEE International Conference on Communications (ICC)*, Beijing, China, May 2008, pp. 2383–2387.
- [132] D. Kim, H. Bae, and C. K. Toh, "Improving TCP-Vegas performance over MANET routing protocols," *IEEE Transactions on Vehicular Technology*, vol. 56, no. 1, pp. 372–377, Jan. 2007.
- [133] J. Sing and B. Soh, "TCP New Vegas revisited," in 13th IEEE International Conference on Networks, vol. 2, Nov. 2005.
- [134] Y.-M. Liu and X.-H. Jiang, "An extended DCCP congestion control in Wireless Sensor Networks," in *IEEE International Workshop on Intelligent Systems and Applications (ISA)*, Wuhan, China, May 2009, pp. 1–4.
- [135] L. Ye and Z. Wang, "A QoS-aware congestion control mechanism for DCCP," in *IEEE Symposium on Computers and Communications* (*ISCC*), Sousse, Tunisia, Jul. 2009, pp. 624–629.
- [136] G. Hasegawa, K. Kurata, and M. Murata, "Analysis and improvement of fairness between TCP Reno and Vegas for deployment of TCP Vegas to the Internet," in *8th IEEE International Conference on Network Protocols (ICNP)*, Osaka, Japan, Nov. 2000, pp. 177–186.
- [137] A. Aijaz, M. Dohler, A. H. Aghvami, V. Friderikos, and M. Frodigh, "Realizing the tactile internet: Haptic communications over next generation 5G cellular networks," *IEEE Wireless Communications*, vol. 24, no. 2, pp. 82–89, April 2017.
- [138] V. Arun and H. Balakrishnan, "Copa: Practical delay-based congestion control for the Internet," in 15th USENIX Symposium on Networked Systems Design and Implementation (NSDI). Renton, Washington, USA: USENIX Association, Apr. 2018, pp. 329–342.
- [139] R. Mittal, V. T. Lam, N. Dukkipati, E. Blem, H. Wassel, M. Ghobadi, A. Vahdat, Y. Wang, D. Wetherall, and D. Zats, "TIMELY: RTT-based congestion control for the datacenter," *ACM Computer Communication Review*, vol. 45, no. 4, pp. 537–550, Aug. 2015.
- [140] A. Zanella, G. Procissi, M. Gerla, and M. Sanadidi, "TCP Westwood: Analytic model and performance evaluation," in *IEEE Global Telecommunications Conference (GLOBECOM)*, vol. 3. San Antonio, Texas, USA: IEEE, 2001, pp. 1703–1707.
- [141] S. Mascolo, L. A. Grieco, R. Ferorelli, P. Camarda, and G. Piscitelli, "Performance evaluation of Westwood+ TCP congestion control," *Performance Evaluation*, vol. 55, no. 1-2, pp. 93–111, Jan. 2004.
- [142] L. Kleinrock, "Power and deterministic rules of thumb for probabilistic problems in computer communications," in *IEEE International Conference on Communications (ICC)*, Boston, Massachusetts, USA, Jun. 1979, pp. 43.1.1–43.1.10.
- [143] J. Jaffe, "Flow control power is nondecentralizable," *IEEE Transactions on Communications*, vol. 29, no. 9, pp. 1301–1306, Sep. 1981.
- [144] M. Hock, R. Bless, and M. Zitterbart, "Experimental evaluation of BBR congestion control," in 25th IEEE International Conference on

*Network Protocols (ICNP)*, Toronto, Ontario, Canada, Oct. 2017, pp. 1–10.

22

- [145] K. Miyazawa, K. Sasaki, N. Oda, and S. Yamaguchi, "Cyclic performance fluctuation of TCP BBR," in 42nd IEEE Annual Computer Software and Applications Conference (COMPSAC), vol. 1, Tokyo, Japan, Jul. 2018, pp. 811–812.
- [146] P. Farrow, "Performance analysis of heterogeneous TCP congestion control environments," in *IEEE/IFIP International Conference on Performance Evaluation and Modeling in Wired and Wireless Networks* (*PEMWN*), Paris, France, Nov. 2017, pp. 1–6.
- [147] Z. Zhong, I. Hamchaoui, R. Khatoun, and A. Serhrouchni, "Performance evaluation of CQIC and TCP BBR in mobile network," in 21st IEEE Conference on Innovation in Clouds, Internet and Networks and Workshops (ICIN), Paris, France, Feb. 2018, pp. 1–5.
- [148] E. Atxutegi, F. Liberal, H. K. Haile, K. Grinnemo, A. Brunstrom, and A. Arvidsson, "On the use of TCP BBR in cellular networks," *IEEE Communications Magazine*, vol. 56, no. 3, pp. 172–179, Mar. 2018.
- [149] F. Li, J. W. Chung, X. Jiang, and M. Claypool, "TCP CUBIC versus BBR on the Highway," in *International Conference on Passive and Active Network Measurement (PAM)*. Cleveland, Ohio, USA: Springer, 2018, pp. 269–280.
- [150] N. Cardwell, Y. Cheng, C. S. Gunn, S. H. Yeganeh, and V. Jacobson, "A quick update on BBR in shallow buffers," in *IETF 100*, Singapore, Nov. 2017. [Online]. Available: https://datatracker.ietf.org/meeting/100/ materials/slides-100-iccrg-a-quick-bbr-update-bbr-in-shallow-buffers
- [151] M. Wang, Y. Cui, X. Wang, S. Xiao, and J. Jiang, "Machine learning for networking: Workflow, advances and opportunities," *IEEE Network*, vol. 32, no. 2, pp. 92–99, Mar. 2018.
- [152] M. Dong, T. Meng, D. Zarchy, E. Arslan, Y. Gilad, B. Godfrey, and M. Schapira, "PCC Vivace: Online-learning congestion control," in 15th USENIX Symposium on Networked Systems Design and Implementation (NSDI), Renton, Washington, USA, Apr. 2018, pp. 343–356.
- [153] Y. Kong, H. Zang, and X. Ma, "Improving TCP congestion control with machine intelligence," in ACM SIGCOMM Workshop on Network Meets AI&ML (NetAI). Budapest, Hungary: ACM, Aug. 2018, pp. 60–66.
- [154] W. Li, F. Zhou, W. Meleis, and K. Chowdhury, "Learning-based and data-driven TCP design for memory-constrained IoT," in *IEEE International Conference on Distributed Computing in Sensor Systems* (DCOSS), Washington, DC, USA, May 2016, pp. 199–205.
- [155] A. P. Silva, K. Obraczka, S. Burleigh, and C. M. Hirata, "Smart congestion control for delay- and disruption tolerant networks," in 13th IEEE International Conference on Sensing, Communication, and Networking (SECON), London, UK, Jun. 2016, pp. 1–9.
- [156] M. Schapira and K. Winstein, "Congestion-Control Throwdown," in 16th ACM Workshop on Hot Topics in Networks (HotNets), Palo Alto, California, USA, Nov. 2017, pp. 122–128.
- [157] M. Zhang, M. Mezzavilla, J. Zhu, S. Rangan, and S. Panwar, "TCP dynamics over mmWave links," in 18th IEEE International Workshop on Signal Processing Advances in Wireless Communications (SPAWC), Sapporo, Japan, July 2017, pp. 1–6.
- [158] T. Azzino, M. Drago, M. Polese, A. Zanella, and M. Zorzi, "X-TCP: A cross layer approach for TCP uplink flows in mmWave networks," in 16th Annual Mediterranean Ad Hoc Networking Workshop (Med-Hoc-Net), Budva, Montenegro, Jun. 2017, pp. 1–6.
- [159] M. Polese, M. Mezzavilla, M. Zhang, J. Zhu, S. Rangan, S. Panwar, and M. Zorzi, "milliProxy: A TCP proxy architecture for 5G mmWave cellular systems," in *51st Asilomar Conference on Signals, Systems, and Computers*, Pacific Grove, California, USA, Oct. 2017, pp. 951–957.
- [160] A. T. Naman, Y. Wang, H. H. Gharakheili, V. Sivaraman, and D. Taubman, "Responsive high throughput congestion control for interactive applications over SDN-enabled networks," *Computer Networks*, vol. 134, pp. 152–166, Apr. 2018.
- [161] H. Shiang and M. van der Schaar, "A quality-centric TCP-friendly congestion control for multimedia transmission," *IEEE Transactions* on Multimedia, vol. 14, no. 3, pp. 896–909, Jun. 2012.
- [162] O. Habachi, H.-P. Shiang, M. van der Schaar, and Y. Hayel, "Online learning based congestion control for adaptive multimedia transmission," *IEEE Transactions on Signal Processing*, vol. 61, no. 6, pp. 1460–1469, Mar. 2013.
- [163] S. M. Aghdam, M. Khansari, H. R. Rabiee, and M. Salehi, "WCCP: A congestion control protocol for wireless multimedia communication in sensor networks," *Ad Hoc Networks*, vol. 13, pp. 516–534, February 2014.
- [164] F. Lu, H. Du, A. Jain, G. M. Voelker, A. C. Snoeren, and A. Terzis, "CQIC: Revisiting cross-layer congestion control for cellular networks," in 16th ACM International Workshop on Mobile Computing

Systems and Applications (HotMobile). Santa Fe, New Mexico, USA: ACM, Feb. 2015, pp. 45–50.

- [165] C. Wilson, H. Ballani, T. Karagiannis, and A. Rowtron, "Better never than late: Meeting deadlines in datacenter networks," in ACM Computer Communication Review, vol. 41, no. 4, Aug. 2011, pp. 50–61.
- [166] D. Drutskoy, E. Keller, and J. Rexford, "Scalable network virtualization in Software-Defined Networks," *IEEE Internet Computing*, vol. 17, no. 2, pp. 20–27, Mar. 2013.
- [167] B. Vamanan, J. Hasan, and T. Vijaykumar, "Deadline-aware Datacenter TCP (D<sup>2</sup>TCP)," ACM Computer Communication Review, vol. 42, no. 4, pp. 115–126, Sep. 2012.
- [168] M. Alizadeh, S. Yang, M. Sharif, S. Katti, N. McKeown, B. Prabhakar, and S. Shenker, "pFabric: Minimal near-optimal datacenter transport," in *ACM Computer Communication Review*, vol. 43, no. 4. ACM, Aug. 2013, pp. 435–446.
- [169] P. X. Gao, A. Narayan, G. Kumar, R. Agarwal, S. Ratnasamy, and S. Shenker, "pHost: Distributed near-optimal datacenter transport over commodity network fabric," in *11th ACM Conference on Emerging Networking Experiments and Technologies (CoNEXT)*. Heidelberg, Germany: ACM, Dec. 2015, pp. 1–12.
- [170] I. Cho, K. Jang, and D. Han, "Credit-scheduled delay-bounded congestion control for datacenters," in ACM Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications (SIGCOMM). Los Angeles, California, USA: ACM, Aug. 2017, pp. 239–252.
- [171] B. Montazeri, Y. Li, M. Alizadeh, and J. Ousterhout, "Homa: A receiver-driven low-latency transport protocol using network priorities," in ACM Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications (SIGCOMM). Budapest, Hungary: ACM, Aug. 2018, pp. 221–235.
- [172] A. Munir, G. Baig, S. M. Irteza, I. A. Qazi, A. X. Liu, and F. R. Dogar, "Friends, not foes: synthesizing existing transport strategies for data center networks," *ACM Computer Communication Review*, vol. 44, no. 4, pp. 491–502, Feb. 2015.
- [173] M. Noormohammadpour and C. S. Raghavendra, "Datacenter traffic control: Understanding techniques and tradeoffs," *IEEE Communications Surveys and Tutorials*, vol. 20, no. 2, pp. 1492–1525, Dec. 2017.
- [174] J. Zhang, F. R. Yu, S. Wang, T. Huang, Z. Liu, and Y. Liu, "Load balancing in data center networks: A survey," *IEEE Communications Surveys and Tutorials*, vol. 20, no. 3, pp. 2324–2352, Jul. 2018.
- [175] T. Hafeez, N. Ahmed, B. Ahmed, and A. W. Malik, "Detection and mitigation of congestion in SDN enabled data center networks: A survey," *IEEE Access*, vol. 6, pp. 1730–1740, 2018.
- [176] A. Ford, C. Raiciu, M. Handley, O. Bonaventure, and C. Paasch, "TCP extensions for multipath operation with multiple addresses," IETF, RFC 6824, Jan. 2013. [Online]. Available: https://rfc-editor.org/ rfc/rfc6824.txt
- [177] O. Mehani, R. Holz, S. Ferlin, and R. Boreli, "An early look at Multipath TCP deployment in the wild," in 6th ACM International Workshop on Hot Topics in Planet-Scale Measurement (HotPlanet), Paris, France, 2015, pp. 7–12.
- [178] C. Raiciu, C. Paasch, S. Barre, A. Ford, M. Honda, F. Duchene, O. Bonaventure, and M. Handley, "How hard can it be? designing and implementing a deployable Multipath TCP," in *9th USENIX Conference on Networked Systems Design and Implementation (NSDI)*, Apr. 2012, pp. 29–29.
- [179] M. Scharf and A. Ford, "Multipath TCP (MPTCP) application interface considerations," IETF, RFC 6897, Mar. 2013. [Online]. Available: https://rfc-editor.org/rfc/rfc6897.txt
- [180] M. Bagnulo, "Threat analysis for TCP extensions for multipath operation with multiple addresses," IETF, RFC 6181, Mar. 2011. [Online]. Available: https://rfc-editor.org/rfc/rfc6181.txt
- [181] S. Barré, C. Paasch, and O. Bonaventure, "Multipath TCP: From theory to practice," in *IEEE/IFIP Networking Conference (IFIP Networking)*, Valencia, Spain, May 2011, pp. 444–457.
- [182] O. Bonaventure, C. Paasch, and G. Detal, "Use cases and operational experience with Multipath TCP," IETF, RFC 8041, Jan. 2017. [Online]. Available: https://rfc-editor.org/rfc/rfc8041.txt
- [183] B. Han, F. Qian, and L. Ji, "When should we surf the mobile web using both WiFi and cellular?" in 5th ACM SIGCOMM Workshop on All Things Cellular: Operations, Applications and Challenges (AllThingsCellular). London, UK: ACM, 2016, pp. 7–12.
- [184] M. Li, A. Lukyanenko, and Y. Cui, "Network coding based multipath TCP," in *IEEE Conference on Computer Communications Workshops* (*INFOCOM WKSHPS*), Orlando, Florida, USA, Mar. 2012, pp. 25–30.
- [185] K. Nguyen, G. P. Villardi, M. G. Kibria, K. Ishizu, F. Kojima, and H. Shinbo, "An enhancement of Multipath TCP performance in lossy

wireless networks," in 41st IEEE Conference on Local Computer Networks Workshops (LCN Workshops), Dubai, UAE, Nov. 2016, pp. 187–191.

- [186] M. Becke, T. Dreibholz, H. Adhari, and E. P. Rathgeb, "On the fairness of transport protocols in a multi-path environment," in *IEEE International Conference on Communications (ICC)*, Ottawa, Ontario, Canada, Jun. 2012, pp. 2666–2672.
- [187] H. Han, S. Shakkottai, C. V. Hollot, R. Srikant, and D. Towsley, "Multipath TCP: a joint congestion control and routing scheme to exploit path diversity in the internet," *IEEE/ACM Transactions on Networking*, vol. 14, no. 6, pp. 1260–1271, Dec. 2006.
- [188] D. Wischik, C. Raiciu, A. Greenhalgh, and M. Handley, "Design, implementation and evaluation of congestion control for Multipath TCP," in 8th USENIX Conference on Networked Systems Design and Implementation (NSDI), Boston, Massachusetts, USA, 2011, pp. 99– 112.
- [189] R. Khalili, N. Gast, M. Popovic, and J.-Y. Le Boudec, "MPTCP is not Pareto-optimal: Performance issues and a possible solution," *IEEE/ACM Transantions on Networking*, vol. 21, no. 5, pp. 1651–1665, Oct. 2013.
- [190] A. Walid, J. Hwang, Q. Peng, and S. Low, "Balanced Linked Adaptation congestion control algorithm for MPTCP," IETF, Working Draft draft-walid-mptcp-congestion-control-00, Jul. 2014. [Online]. Available: http://www.ietf.org/internet-drafts/ draft-walid-mptcp-congestion-control-00.txt
- [191] M. Zhu, L. Wang, Z. Qin, N. Ding, J. Fang, T. Liu, and Q. Cui, "BE-LIA: bandwidth estimate-based link increase algorithm for MPTCP," *IET Networks*, vol. 6, no. 5, pp. 94–101, September 2017.
- [192] Z. Xu, J. Tang, C. Yin, Y. Wang, and G. Xue, "Experience-Driven Congestion Control: When Multi-Path TCP Meets Deep Reinforcement Learning," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 6, pp. 1325–1336, June 2019.
- [193] J. R. Iyengar, P. D. Amer, and R. Stewart, "Concurrent multipath transfer using SCTP multihoming over independent end-to-end paths," *IEEE/ACM Transactions on Networking*, vol. 14, no. 5, pp. 951–964, Oct. 2006.
- [194] M. Becke, H. Adhari, E. P. Rathgeb, F. Fa, X. Yang, and X. Zhou, "Comparison of Multipath TCP and CMT-SCTP based on intercontinental measurements," in *IEEE Global Communications Conference* (*GLOBECOM*). Atlanta, Georgia, USA: IEEE, 2013, pp. 1360–1366.
- [195] C.-M. Huang, Y.-C. Chen, and S.-Y. Lin, "Packet scheduling and congestion control schemes for Multipath Datagram Congestion Control Protocol," *The Computer Journal*, vol. 58, no. 2, pp. 188–203, Feb. 2015.
- [196] B. Y. L. Kimura, D. C. S. F. Lima, and A. A. F. Loureiro, "Alternative scheduling decisions for Multipath TCP," *IEEE Communications Letters*, vol. 21, no. 11, pp. 2412–2415, Nov. 2017.
- [197] D. Ni, K. Xue, P. Hong, and S. Shen, "Fine-grained forward prediction based dynamic packet scheduling mechanism for multipath TCP in lossy networks," in 23rd IEEE International Conference on Computer Communication and Networks (ICCCN), Shanghai, China, 2014, pp. 1–7.
- [198] K. W. Choi, Y. S. Cho, J. W. Lee, S. M. Cho, J. Choi et al., "Optimal load balancing scheduler for MPTCP-based bandwidth aggregation in heterogeneous wireless environments," *Computer Communications*, vol. 112, pp. 116–130, Nov. 2017.
- [199] Y.-s. Lim, E. M. Nahum, D. Towsley, and R. J. Gibbens, "ECF: An MPTCP Path Scheduler to Manage Heterogeneous Paths," in 13th ACM International Conference on Emerging Networking Experiments and Technologies (CoNEXT), Incheon, South Korea, Dec. 2017, pp. 147– 159.
- [200] Y. E. Guo, A. Nikravesh, Z. M. Mao, F. Qian, and S. Sen, "Demo: DEMS: DEcoupled Multipath Scheduler for accelerating multipath transport," in 23rd ACM International Conference on Mobile Computing and Networking (MobiCom), Snowbird, Utah, USA, Oct. 2017, pp. 477–479.
- [201] H. Shi, Y. Cui, X. Wang, Y. Hu, M. Dai, F. Wang, and K. Zheng, "STMS: Improving MPTCP throughput under Heterogeneous Networks," in USENIX Annual Technical Conference (USENIX ATC), Boston, Massachusetts, USA, Jul. 2018, pp. 719–730.
- [202] N. Kuhn, E. Lochin, A. Mifdaoui, G. Sarwar, O. Mehani, and R. Boreli, "DAPS: Intelligent delay-aware packet scheduling for multipath transport," in *IEEE International Conference on Communications (ICC)*. IEEE, 2014, pp. 1222–1227.
- [203] S. Ferlin, O. Alay, O. Mehani, and R. Boreli, "BLEST: Blocking estimation-based MPTCP scheduler for heterogeneous networks," in

IEEE/IFIP Networking Conference (IFIP Networking), Vienna, Austria, May 2016, pp. 431–439.

- [204] K. Fahmi, D. Leith, S. Kucera, and H. Claussen, "Low Delay Scheduling of Objects Over Multiple Wireless Paths," arXiv preprint arXiv:1808.02418, Aug. 2018. [Online]. Available: https: //arxiv.org/abs/1808.02418
- [205] T. Shreedhar, N. Mohan, S. Kaul, J. Kangasharju et al., "QAware: a cross-layer approach to MPTCP scheduling," in *IEEE/IFIP Networking Conference (IFIP Networking)*, Zürich, Switzerland, May 2018.
- [206] M. Xu, Y. Cao, and E. Dong, "Delay-based congestion control for MPTCP," IETF, Working Draft: draft-xu-mptcp-congestion-control-05, Jan. 2017. [Online]. Available: https://datatracker.ietf.org/doc/html/ draft-xu-mptcp-congestion-control-05
- [207] T.-A. Le, "Improving the performance of multipath congestion control over wireless networks," in *International Conference on Advanced Technologies for Communications (ATC)*. Vietri sul Mare, Italy: IEEE, Dec. 2013, pp. 60–65.
- [208] T. Shreedhar, N. Mohan, S. K. Kaul, and J. Kangasharju, "More than the sum of its parts: Exploiting cross-layer and joint-flow information in MPTCP," arXiv preprint arXiv:1711.07565, Nov. 2017. [Online]. Available: https://arxiv.org/abs/1711.07565
- [209] K. Noda, Y. Ito, and Y. Muraki, "Study on congestion control of multipath TCP based on web-QoE under heterogeneous environment," in 6th IEEE Global Conference on Consumer Electronics (GCCE), Nagoya, Japan, Oct. 2017, pp. 1–3.
- [210] Y. Muraki and Y. Ito, "Study on effect of congestion control of multipath TCP on Web-QoE," in 4th IEEE Global Conference on Consumer Electronics (GCCE), Osaka, Japan, Oct. 2015, pp. 52–53.
- [211] Y.-C. Chen, Y.-s. Lim, R. J. Gibbens, E. M. Nahum, R. Khalili, and D. Towsley, "A measurement-based study of Multipath TCP performance over wireless networks," in ACM Conference on Internet Measurement, 2013, pp. 455–468.
- [212] M. Li, A. Lukyanenko, S. Tarkoma, Y. Cui, and A. Ylä-Jääski, "Tolerating path heterogeneity in Multipath TCP with bounded receive buffers," in ACM International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS), Pittsburgh, PA, USA, Jun. 2013, pp. 375–376.
- [213] A. Garcia-Saavedra, M. Karzand, and D. J. Leith, "Low Delay Random Linear Coding and Scheduling Over Multiple Interfaces," *IEEE Transactions on Mobile Computing*, vol. 16, no. 11, pp. 3100–3114, Nov. 2017.
- [214] P. Ageneau, N. Boukhatem, and M. Gerla, "Practical random linear coding for MultiPath TCP: MPC-TCP," in 24th IEEE International Conference on Telecommunications (ICT), Limassol, Cyprus, May 2017, pp. 1–6.
- [215] Y. Cui, L. Wang, X. Wang, H. Wang, and Y. Wang, "FMTCP: A fountain code-based Multipath Transmission Control Protocol," *IEEE/ACM Transactions on Networking*, vol. 23, no. 2, pp. 465–478, Apr. 2015.
- [216] S. Ferlin, S. Kucera, H. Claussen, and O. Alay, "MPTCP meets FEC: Supporting latency-sensitive applications over Heterogeneous Networks," *IEEE/ACM Transactions on Networking*, pp. 1–14, Oct. 2018.
- [217] S. Ferlin, T. Dreibholz, and Ö. Alay, "Multi-path transport over heterogeneous wireless networks: Does it really pay off?" in *IEEE Global Communications Conference (GLOBECOM)*. Austin, Texas, USA: IEEE, Dec. 2014, pp. 4807–4813.
- [218] M. S. Kim, I. H. Jung, K. M. Han, J. Y. Lee, and B. C. Kim, "Performance enhancement of MPTCP having a bufferbloat path using retransmission of HoL blocking packets," in *IEEE International Conference on Electronics, Information, and Communication (ICEIC)*, Honolulu, Hawaii, USA, Jan. 2018, pp. 1–2.
- [219] J. Hwang and J. Yoo, "Packet scheduling for Multipath TCP," in 7th International Conference on Ubiquitous and Future Networks, Jul. 2015, pp. 177–179.
- [220] E. Dong, M. Xu, X. Fu, and Y. Cao, "LAMPS: A Loss Aware Scheduler for Multipath TCP over Highly Lossy Networks," in 42nd IEEE Conference on Local Computer Networks (LCN), Singapore, Oct. 2017, pp. 1–9.
- [221] C. Paasch, S. Ferlin, O. Alay, and O. Bonaventure, "Experimental evaluation of multipath TCP schedulers," in ACM SIGCOMM Workshop on Capacity Sharing (CSWS). Chicago, Illinois, USA: ACM, Aug. 2014, pp. 27–32.
- [222] C. Raiciu, C. Pluntke, S. Barre, A. Greenhalgh, D. Wischik, and M. Handley, "Data center networking with multipath TCP," in 9th ACM Workshop on Hot Topics in Networks (HotNets), Monterey, California, USA, Oct. 2010, p. 10.

- [223] S. Tariq and M. Bassiouni, "Performance evaluation of MPTCP over optical burst switching in data centers," in *IEEE International Telecommunications Symposium (ITS)*, São Paulo, Brazil, Aug. 2014, pp. 1–5.
- [224] M. Li, A. Lukyanenko, S. Tarkoma, and A. Ylä-Jääski, "MPTCP incast in data center networks," *China Communications*, vol. 11, no. 4, pp. 25–37, Apr. 2014.
- [225] J. Hwang, A. Walid, and J. Yoo, "Fast coupled retransmission for multipath TCP in data center networks," *IEEE Systems Journal*, vol. 12, no. 1, pp. 1056–1059, Mar. 2018.
- [226] S. Zannettou, M. Sirivianos, and F. Papadopoulos, "Exploiting path diversity in datacenters using MPTCP-aware SDN," in 21st IEEE Symposium on Computers and Communication (ISCC), Messina, Italy, Jun. 2016, pp. 539–546.
- [227] J. Duan, Z. Wang, and C. Wu, "Responsive multipath TCP in SDNbased datacenters," in *IEEE International Conference on Communications (ICC)*, London, UK, Jun. 2015, pp. 5296–5301.
- [228] T. Dreibholz, M. Becke, H. Adhari, and E. P. Rathgeb, "On the impact of congestion control for concurrent multipath transfer on the transport layer," in *11th IEEE International Conference on Telecommunications* (*ConTEL*), Graz, Austria, Jun. 2011, pp. 397–404.
- [229] Q. De Coninck and O. Bonaventure, "Multipath QUIC: Design and evaluation," in 13th ACM International Conference on Emerging Networking EXperiments and Technologies (CoNEXT), Incheon, South Korea, Dec. 2017, pp. 160–166.
- [230] M. Amend, A. Brunstrom, A. Kassler, and V. Rakocevic, "IP compatible multipath framework for heterogeneous access networks," IETF, Working Draft draft-amend-tsvwg-multipath-framework-mpdccp-00, Mar. 2019. [Online]. Available: http://www.ietf.org/internet-drafts/ draft-amend-tsvwg-multipath-framework-mpdccp-00.txt

![](_page_23_Picture_30.jpeg)

Michele Polese [S'17] received his B.Sc. in Information Engineering in 2014 and M.Sc. in Telecommunication Engineering in 2016 from the University of Padova, Italy. Since October 2016 he has been a Ph.D. student at the Department of Information Engineering of the University of Padova, under the supervision of Prof. Michele Zorzi. He visited New York University (NYU), NY in 2017, AT&T Labs in Bedminster, NJ in 2018, and Northeastern University, Boston, MA in 2019. He is collaborating with several academic and industrial research partners,

including Intel, InterDigital, NYU, AT&T Labs, University of Aalborg, King's College, Northeastern University and NIST. He was awarded with the Best Journal Paper Award of the IEEE ComSoc Technical Committee on Communications Systems Integration and Modeling (CSIM) 2019 and the Best Paper Award at WNS3 2019. His research interests are in the analysis and development of protocols and architectures for the next generation of cellular networks (5G), in particular for millimeter-wave communication, and in the performance evaluation of complex networks.

![](_page_23_Picture_33.jpeg)

Federico Chiariotti [S'15 M'19] is a post-doctoral researcher at the University of Padova, in Italy, where he received his Ph.D. in information engineering in 2019. He received the bachelor's and master's degrees in telecommunication engineering (both cum laude) from the University of Padova, in 2013 and 2015, respectively. In 2017 and 2018, he was a Research Intern with Nokia Bell Labs, Dublin. He has authored over 20 published papers on wireless networks and the use of artificial intelligence techniques to improve their performance.

He was a recipient of the Best Paper Award at the Workshop on ns-3 in 2019 and the Best Student Paper Award at the International Astronautical Congress in 2015. His current research interests include network applications of machine learning, transport layer protocols, Smart Cities, bike sharing system optimization, and adaptive video streaming.

25

![](_page_24_Picture_1.jpeg)

**Elia Bonetto** Elia Bonetto received his B.Sc. degree in computer engineering from the University of Padova, Padova, Italy, in 2017, where he is currently pursuing the M.Sc. degree in ICT for Internet and multimedia. His main interests are in computer vision and robotics besides computer networks and network protocols.

![](_page_24_Picture_3.jpeg)

Filippo Rigotto received his B.Sc. degree in computer engineering from the University of Padova, Italy, in 2017, where he is currently pursuing the M.Sc. degree in ICT for Internet and multimedia. His present interests are in computer networks and network protocols, machine learning, computer vision and robotics areas.

![](_page_24_Picture_5.jpeg)

Andrea Zanella [S'98-M'01-SM'13] is Associate Professor at the University of Padova, in Italy. He received the master degree in Computer Engineering and the Ph.D. degree in Electronic and Telecommunications Engineering from the same University, in 1998 and 2000, respectively. He has (co)authored more than 130 papers, five books chapters and three international patents in multiple subjects related to wireless networking and Internet of Things, Vehicular Networks, cognitive networks, and microfluidic

networking. He serves as Technical Area Editor for the IEEE INTERNET OF THINGS JOURNAL, and as Associate Editor for the IEEE COMMUNICATIONS SURVEYS & TUTORIALS, and the IEEE TRANSACTIONS ON COGNITIVE COMMUNICATIONS AND NETWORKING.

![](_page_24_Picture_8.jpeg)

Michele Zorzi [F'07] received his Laurea and PhD degrees in electrical engineering from the University of Padova in 1990 and 1994, respectively. During academic year 1992-1993 he was on leave at the University of California at San Diego (UCSD). In 1993 he joined the faculty of the Dipartimento di Elettronica e Informazione, Politecnico di Milano, Italy. After spending three years with the Center for Wireless Communications at UCSD, in 1998 he joined the School of Engineering of the University of Ferrara, Italy, where he became a professor in

2000. Since November 2003 he has been on the faculty of the Information Engineering Department at the University of Padova. His present research interests include performance evaluation in mobile communications systems, WSN and Internet of Things, cognitive communications and networking, 5G mmWave cellular systems, vehicular networks, and underwater communications and networks. He is the recipient of several awards from the IEEE Communications Society, including the Best Tutorial Paper Award (2008, 2019), the Education Award (2016), and the Stephen O. Rice Best Paper Award (2018). He was the Editor in Chief of the IEEE Transactions on Cognitive Communications from 2003 to 2005 and of the IEEE Transactions on Communications from 2008 to 2011. He served as a Member-at-Large of the Board of Governors of the IEEE Communications Society from 2009 to 2011, and as its Director of Education in 2014-15.