# Evaluating the Impact of Packet Scheduling and Congestion Control Algorithms on MPTCP Performance over Heterogeneous Networks

DIMITRIOS DIMOPOULOS, Lenovo Research, Greece and Department of Informatics and Telecommunications, National and Kapodistrian University of Athens, Greece

APOSTOLIS K. SALKINTZIS, Lenovo Research, Greece

DIMITRIS TSOLKAS, Department of Informatics and Telecommunications, National and Kapodistrian University of Athens, Greece

NIKOS PASSAS, Department of Informatics and Telecommunications, National and Kapodistrian University of Athens, Greece

LAZAROS MERAKOS, Department of Informatics and Telecommunications, National and Kapodistrian University of Athens, Greece

Modern mobile and stationary devices are equipped with multiple network interfaces aiming to provide wireless and wireline connectivity either in a local LAN or the Internet. Multipath TCP (MPTCP) protocol has been developed on top of legacy TCP to allow the simultaneous use of multiple network paths in the communication route between two end-systems. Although the combination of multiple paths is beneficial in case of links with similar network characteristics, MPTCP performance is challenged as heterogeneity among the used paths increases. This work provides an overview of the MPTCP protocol operation, analyzes the state-of-art packet scheduling and congestion control algorithms available in literature, and examines the impact of the various algorithm combinations on MPTCP performance, by conducting an extensive experimental evaluation under diverse path-heterogeneity conditions.

Additional Key Words and Phrases: MPTCP, packet scheduling, congestion control, performance evaluation

## 1 Introduction

Modern digital devices, from laptops and cell phones to large data centers, are usually equipped with multiple network interfaces (e.g., Ethernet, WiFi, LTE/5G, optical, Bluetooth, etc.). Multipath TCP (MPTCP) protocol [36], which builds upon and extends standard TCP, emerged from the necessity to provide data communication services over multiple paths, by jointly using the available network interfaces of a multi-homed device. Although new protocols arise (e.g, *QUIC* [26, 52] along with its multipath extension *MPQUIC* [65]), TCP is still the predominant protocol for critical applications, as well as the default standard in all operating systems. MPTCP enables the exchange of data traffic over multiple disjoint or common-bottleneck links allowing the communicating hosts to reach specific application goals; be it higher data rate, lower flow completion time, or resilience in case of link failure. To this end, MPTCP provides a clear benefit to the application between two communicating systems. However, the simultaneous use of paths with inherently different network characteristics poses a challenge on MPTCP's performance. In fact, network links are affected either by the constantly fluctuating demand in shared network resources (e.g., in the Internet) which often leads to network congestion, increased intermediate buffer queues and packet drops, or by the nature of the physical

medium (e.g., the 5G wireless interface) which is susceptible to random packet loss. The presence of path heterogeneity, in conjunction with MPTCP's adherence to TCP principles for reliable and in-order bytestream delivery, may lead to head-of-line (HoL) blocking at the receiver, stalling the connection and degrading the overall performance. In an attempt to tackle the functional and performance issues arising by the presence of heterogeneous paths, research efforts focus on providing enhanced congestion control and optimal packet scheduling algorithms that increase the overall performance and are tailored to application demands. In this direction, we hereby provide a comprehensible guide on MPTCP protocol structure and functionality, analyze the theoretical aspects behind the operation of the state-of-art packet scheduling and congestion control algorithms, construct and present the detailed algorithmic steps in pseudo-code, and finally, perform an extensive performance assessment of the algorithms' combined operation under homogeneous and asymmetric path conditions. The remainder of the paper is structured as follows. Section 2 provides the necessary background on MPTCP protocol aspects and analyzes the best to date packet scheduling and congestion control algorithms. Going beyond theory, an extensive experimental evaluation has been conducted, providing insight on how MPTCP performs under different packet scheduling and congestion control mechanism combinations. To this end, Section 3 describes the experimentation methodology, while Section 4 presents the produced results. Section 5 discusses the performance results and provides insight on the best-performing congestion control and packet scheduling schemes, while Section 6 is dedicated to the corresponding works available in literature. Finally, Section 7 concludes the manuscript.

## 2  Background

This section is dedicated in describing the fundamental aspects of MPTCP, such as the protocol structure, the architectural details, as well as foundational knowledge behind the most salient modules comprising and dictating the operation of MPTCP; the Path Manager, the Packet Scheduler, and the Congestion Controller. Afterwards, the main logic, as well as the algorithmic steps of the state-of-art packet scheduling and congestion control algorithms are presented.

### 2.1  MPTCP Theory

*2.1.1  General principles.* MPTCP has been designed on top of TCP, and as such, shares and extends much of its principles and structural details. The intention behind its development was the design of a protocol which provides the traditional single-path TCP transport layer services, such as connection-oriented, reliable, and in-order data delivery, while being able to simultaneously use multiple paths for the transmission of data segments. MPTCP is thus extending traditional TCP, and its operation is determined by specific attributes present within TCP header options; for instance, the TCP option-kind number *30* has been reserved by IANA to indicate whether the TCP header options carry additional MPTCP related fields. While advancing the legacy TCP functionality, it was also fundamental for MPTCP to be able to fall back to TCP, if required. Thus, the new protocol should not incur any modifications to the overlay application which, as with traditional TCP, can establish a MPTCP connection through a single socket. Despite the use of multiple paths in lower layers, each socket supports a single connection between the application and the transport layer. Applications may continue using the existing TCP socket API for MPTCP connections; in that case MPTCP specific parameters are configured in the operating system. However, a new MPTCP socket is also available allowing applications to control MPTCP behavior via socket parameters.

*2.1.2  Subflows.* The MPTCP's foundational principle is the ability to route application traffic over multiple paths, and yet enable seamless communication between two end-systems, as if it were for a single-path TCP connection. The

underlying paths carrying data packets, can be either distinct end-to-end physical links (i.e., "disjoint" paths), or paths which at some point share a common bottleneck link (e.g., "partially disjoint"). As with standard TCP, different paths are uniquely identified by a 4-tuple (source & destination address/port pair). Each data flow carried on an individual path which belongs to the same MPTCP connection, is called *subflow*. Each subflow resembles a standard single-path TCP connection; subflow connection establishment and termination follows the typical single-path TCP procedures. Consequently, a MPTCP connection can comprise one or more subflows, ultimately enabling application communication between two end-hosts over multiple paths. From application's perspective, there is only a single socket bind and accept call to the MPTCP connection. Each of the communicating multipath-capable hosts generates a token, based on the keys exchanged during the initial MPTCP connection initiation procedure, which uniquely identifies each MPTCP connection. Then, each time a new subflow needs to be created and be associated to an existing MPTCP connection, each host binds the 5-tuple of the TCP subflow to the local token of the connection. This allows any port pairs to be used for a connection, eliminating in essence the identification of a subflow based on the 4-tuple, and relying merely on the token, instead [36].

*2.1.3 Multipath-capable hosts.* The establishment of a multipath TCP connection requires that both communicating hosts have multipath capability. MPTCP is already supported by operating systems which base their function on recent versions of the linux kernel [38]. End-hosts negotiate their multipath capabilities (*MP_CAPABLE*) over a typical 3-way handshake, before being able to establish a MPTCP connection. If one of the hosts is not multipath-capable, or in case middleboxes intervene altering either the TCP header options and removing the MP_CAPABLE field or the payload itself, or in case of MPTCP version mismatch, the connection falls back to regular TCP. This allows the establishment of a standard TCP connection, absent any manual intervention.

*2.1.4 Meta- vs. Subflow-level.* In case of MPTCP, transport layer is logically segregated into two parts; the MPTCP "Meta-level" (or "Connection-level"), and the MPTCP "subflow-level". The MPTCP meta-level is responsible for maintaining a single connection point to the application layer, thus MPTCP is still perceived as a single-path TCP connection from application's point of view. The meta-level has all those mechanisms required to allow the presence of a single connection to the upper layer, yet able to distribute (or reassemble) application bytestream across (or from) multiple subflows. At subflow-level, each subflow is essentially handled as an individual single-path TCP connection between end-hosts.

*2.1.5 Protocol Structure.* Building upon and extending traditional TCP, MPTCP protocol operations are defined within TCP header "Options" field. There are different MPTCP fields set within TCP header "Options" field depending on the various MPTCP operations; i.e., *MP_CAPABLE* option is used during MPTCP connection initiation, the *MP_JOIN* option is used to associate additional subflows with an existing MPTCP connection, the *ADD_ADDR* option denotes the address advertisement operation which announces the availability of new interfaces at end-hosts, the Data Sequence Signal (*DSS*) option carries the Data Sequence Mapping, etc. A general overview of the MPTCP protocol stack as well as the respective protocol structure, is depicted in Figures 1 and 2, respectively.

*2.1.6 Path Manager.* As mentioned earlier, each MPTCP subflow is uniquely identified by a 4-tuple (source & destination address/port pair). The MPTCP *Path Manager* module is responsible for advertising the availability or elimination of an IP addresses/ port pair of a multipath-capable host, for establishing and removing subflows between two end-systems, as well as for maintaining the MPTCP connection state. The MPTCP connection state includes details of the active and standby subflows constituting each MPTCP connection, and maps each subflow's 4-tuple to an *Address ID*; the latter is

```
                                              +------------+----------+
                                              | TCP header |  payload |
                                              +------------+----------+

 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|          Source Port          |       Destination Port        |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                        Sequence Number                        |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                     Acknowledgment Number                     |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
| Data  |           |C|E|U|A|P|R|S|F|                            |
| Offset| Rsrvd     |W|C|R|C|S|S|Y|I|            Window          |
|       |           |R|E|G|K|H|T|N|N|                            |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|           Checksum            |         Urgent Pointer        |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                            [Options]                          |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

```
 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|     Kind      |     Length    |Subtype|                       |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+                       |
|                        Subtype-specific data                  |
|                         (variable length)                     |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

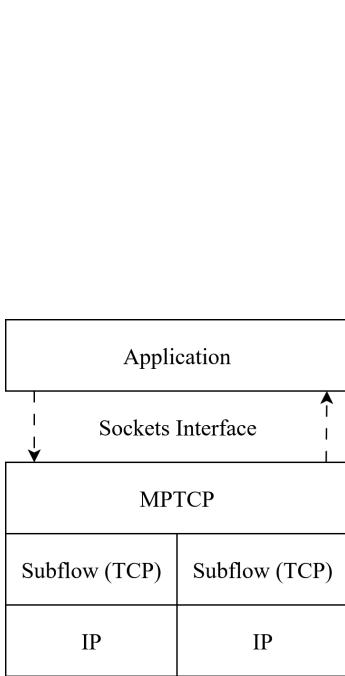| Application |
|---|
| Sockets Interface |
| MPTCP |
| Subflow (TCP) | Subflow (TCP) |
| IP | IP |

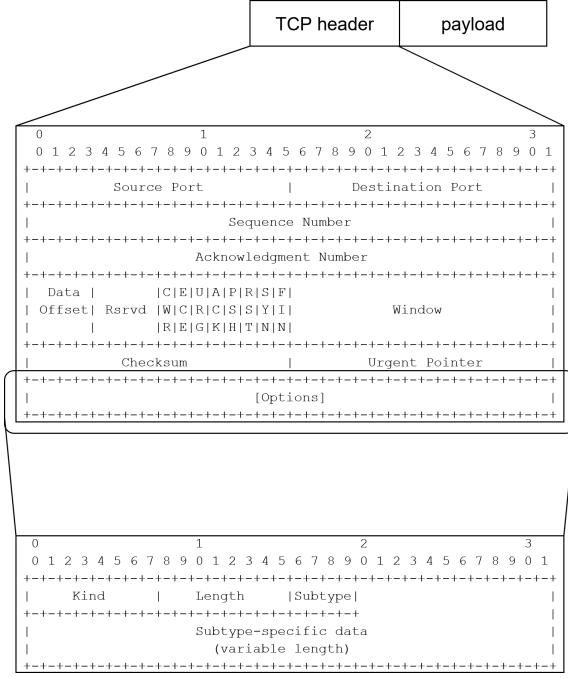Fig. 1. MPTCP protocol stack (Source: IETF RFCs 8684 [36], 6897 [83])

Fig. 2. MPTCP protocol header structure - MPTCP header is included into TCP header "Options" field

generated and announced by a source host during MP_JOIN or ADD_ADDR operation to allow the receiver to refer to a source IP address/port pair implicitly, in case NAT protocol or other middleboxes conceal the actual source details.

*2.1.7  Packet Scheduler.* The *Packet Scheduling* functionality [51] at the transmitting end is responsible for receiving the application bytestream, dividing application data into transport layer segments with appropriate connection-level sequence numbers, and transmitting the segments onto the available subflows according to the scheduler's own policy. The availability of the subflows is exposed to the scheduler by the path manager. At the receiving side, a corresponding functionality is in charge of appropriately ordering the received segments based on the connection-level sequence numbers, and delivering them in a bytestream fashion to the receiving application.

*2.1.8  Flow control.* MPTCP inherits and extends standard TCP mechanisms. *Flow control* guarantees that the transmitted application data will reach the receiving end in-order and flawlessly, while it also allows the receiver to regulate the sender's transmitting rate based on the receiving application processing capabilities; the transmitting rate is essentially controlled via the *receive window (rwnd)* which is announced by the receiver to the sending host at regular intervals (i.e., with each cumulative ACK in TCP). The typical TCP flow control is maintained at subflow-level, where each individual subflow maintains its own *sequence* and *acknowledgment* numbers. However, the existing TCP flow control mechanism is not sufficient since although received packets may end up reassembled in correct order at subflow-level, they may not necessarily be appropriately ordered at MPTCP connection-level due to heterogeneous path conditions or retransmissions. For this reason, besides the subflow-specific sequence and acknowledgment numbers, flow control mechanism is extended at connection-level. *Data Sequence Number (DSN)* and *Data ACK* enable the appropriate *Data*

*Sequence Mapping* from the subflow sequence space to the data sequence space [36]. This extended functionality allows
for a cumulative data tracking at connection-level and ensures in-order delivery to the application.

*2.1.9 Congestion Control - the Reno approach.* The *Congestion control* mechanism has been developed in the context
of TCP [2, 31] to regulate the transmission rate of each sender in an attempt to protect networks, and subsequently
the Internet, from *congestion collapse*. Since then, various congestion control algorithms have emerged, with CUBIC
being the one most broadly utilized by Internet hosts, at the time of this publication. The remainder of this subsection
provides the foundational aspects of congestion control as defined by the legacy TCP-Reno [2]; CUBIC [89] is analyzed
in a dedicated subsection below 2.3.1. TCP's congestion control mechanism comprises the following algorithms (or
"phases"): *slow start*, *congestion avoidance*, *fast retransmit*, and *fast recovery*; the first two are responsible for regulating
the amount of data each TCP sender can inject into the network at any given moment, while the last two determine
the way congestion control reacts to potential loss. Three important state variables influence the function of these
algorithms: i) the congestion window (cwnd), which is essentially a sender-side limit on the amount of data that can be
transmitted and, consequently be in-flight, at any moment, ii) the receive window (rwnd), which is part of flow control
mechanism and constitutes a receiver-side limit on the amount of data the receiver can accept at any given time, and
iii) the slow start threshold (ssthresh) which dictates whether congestion control operates according to slow start or
congestion avoidance algorithm. At any given moment, TCP must send no more than $min(cwnd, rwnd)$ data beyond
the highest acknowledged sequence number; this formula guarantees that the amount of data transmitted at any instant
neither incurs network congestion nor overloads the receive buffer. At any given time, the amount of bytes the sender is
allowed to send next, is defined according to below formula:

$$\#\_of\_bytes\_next = min(cwnd, rwnd) - FlightSize \tag{1}$$

where *FlightSize* is the amount of data in-flight, or in other words, the amount of outstanding data in the network.

It is also recommended that the initial value of ssthresh be set arbitrarily high to allow network conditions, rather
than host limitations, dictate the sending rate. Nevertheless, sshthresh is adjusted to lower values in case congestion is
detected. The selection of the algorithm to be employed by the congestion control mechanism, is based on the relation
between the slow start threshold and the congestion window:

- if $cwnd < ssthresh$, then *slow start* algorithm is employed
- if $cwnd > ssthresh$, then *congestion avoidance* algorithm is employed
- if $cwnd == ssthresh$, then either *slow start* or *congestion avoidance* algorithm is selected

*Slow start.* Slow start algorithm [2] is executed at the beginning of a data transfer after a TCP connection has been
successfully established, or when *retransmission timeout (RTO)* timer elapses, indicating a loss. Slow start algorithm
essentially allows sending-host's TCP to slowly probe the network to estimate network capacity, while avoiding network
congestion induced by large data bursts. In slow start, the initial value of the congestion window immediately after the
completion of the 3-way handshake is:

- if $SMSS > 2190\ bytes$:
    $IW = 2 * SMSS\ bytes$ and NOT more than 2 segments
- if $(SMSS > 1095 bytes)\ \&\&\ (SMSS <= 2190 bytes)$:
    $IW = 3 * SMSS\ bytes$ and NOT more than 3 segments
- if $SMSS <= 1095\ bytes$:

$$IW = 4 * SMSS \; bytes \text{ and NOT more than 4 segments}$$

where *SMSS* is the *Sender Maximum Segment Size*, and *IW* denotes the *Initial Congestion Window*. During slow start, TCP increments the congestion window by at most *SMSS* bytes for each ACK received that cumulatively acknowledges new data:

$$cwnd+ = SMSS \tag{2}$$

A more conservative, yet secure against "ACK Division" approach [82], recommends that TCP increment congestion window according to the actual amount of bytes corresponding to new data, that are successfully acknowledged by the receiver; this transforms 2 to a more precise formula:

$$cwnd+ = min(N, SMSS) \tag{3}$$

where $N$ is the number of previously unacknowledged bytes acknowledged in the incoming ACK. Incrementing the congestion window by at most one SMSS for every received ACK acknowledging new data, leads to an exponential growth of the congestion window, until the point it reaches slow-start threshold (*ssthresh*); from that point onward, *congestion avoidance* algorithm controls data transmission, provided that no loss is encountered.

*Congestion avoidance.* [2] When $cwnd \geq sshthresh$, congestion control enters the congestion avoidance phase. During congestion avoidance, it is recommended that congestion window be incremented proportionally to the number of bytes acknowledged for new data, but no more than *SMSS* bytes per RTT, as defined in equation 4:

$$cwnd + = \; SMSS * SMSS/cwnd \tag{4}$$

Equation 4 provides a byte-unit calculation of the congestion window during CA phase. An alternative to this, is the full-size segment equivalent provided through 3. However, of vital importance here is the fact that the increment takes place per RTT, and not per ACK as is the case in slow start mode.

*Loss detection due to RTO expiration.* [2] When a TCP sender presumes that a segment is lost, by means of RTO expiration, and the segment has not yet been resent for the same reason (i.e., RTO expiration), *ssthresh* is adjusted based on the below equation:

$$ssthresh = \; max\left(\frac{FlightSize}{2}, \; 2 * SMSS\right) \tag{5}$$

If the lost segment has previously been sent as part of retransmission timer expiration, *ssthresh* value remains unmodified. In addition, upon *RTO* expiration, congestion window should be set as follows:

$$cwnd = \; 1 * SMSS \tag{6}$$

At this point, the TCP sender is expected to retransmit the dropped segment and then, using the slow start algorithm, to increase the congestion window from one segment to the new value of *ssthresh* 5. Once *sshthresh* is reached, the algorithm switches to *congestion avoidance* phase [73].

*Fast Retransmit & Fast Recovery.* [2] Fast retransmit algorithm exploits ACKs to detect and repair loss. The receiver should respond with a duplicate-ACK when an out-of-order segment arrives, and should continue responding with duplicate ACKs for every new segment arriving, other than the one still missing. Duplicate-ACKs can in principle be implying a couple of conditions:

   i. the segment has been dropped in transit to the destination

    ii. the network might have re-ordered the segments

    iii. the network replicates ACKs or data segments (i.e., triggered by intermediate devices in the network path between two end nodes)

Upon reception of the missing segment, the receiver is expected to immediately respond with an ACK, indicating that the received segment fills in the gap in the sequence space. Upon reception of the first and second duplicate ACKs, the sender (unless SACK mechanism is enabled[1]) is expected to send a segment of new data, provided that: i) $rwnd$ allows, ii) $FlightSize + min(N, SMSS) \leq cwnd + 2 * SMSS$, and iii) new data is available for transmission. However, $cwnd$ does not yet change, avoiding to reflect the new segments [33].

The reception of three duplicate ACKs denotes that a particular segment has potentially been lost. In that case, the sender sets $ssthresh$ according to equation 5, retransmits immediately the missing segment without waiting for the RTO to expire, and sets $cwnd$ according to the below formula:

$$cwnd = ssthresh + 3 * SMSS \tag{7}$$

This immediate response to the detected loss is termed *fast retransmit*. Equation 7 indicates that the congestion window is artificially inflated by the number of segments (i.e., three) that have left the network and which the receiver has buffered. Notable here is the fact that congestion control does not switch to the slow start phase; this is reasonable, since the reception of duplicate ACKs, if not excessive for different segments, indicates that the network can still deliver segments despite the perceived congestion. As a result, the reception of duplicate ACKs is considered a *transient* congestion, contrary to the expiration of RTO which indicates severe congestion and requests for more drastic measures, enforced by entering the slow start phase.

Following the *fast retransmit*, congestion control enters the *fast recovery* phase, until a normal (i.e., non-duplicate) ACK arrives, indicating the recovery of the segment which was supposed to be missing. In *fast recovery* phase, for each additional duplicate ACK received after the third, congestion window is incremented by one SMSS, that is $cwnd+ = SMSS$, per dup-ACK receipt. This increase artificially inflates the congestion window in order to reflect the segment that has left the network. Based on this increase, the sender can send $1 * SMSS$ bytes of previously unsent data, given that the normal requisites hold true (e.g., new value of $rwnd$ and $cwnd$ allow, and new data is available for transmission). Upon reception of an ACK that acknowledges new data, beyond the sequence range of the missing segment, loss is deemed repaired and TCP sets $cwnd = ssthresh$; beyond that point, *congestion avoidance* algorithm takes over to regulate $cwnd$ increase.

*Final remarks.* The aforementioned congestion control operation is also known as Additive Increase - Multiplicative Decrease (AIMD). The additive increase describes the behavior of the congestion control algorithm within the congestion avoidance phase, where the congestion window increases linearly by at most $1 * SMSS$ bytes per RTT. The multiplicative decrease is referred to loss detection by means of three duplicate ACKs' receipt, where the congestion window is halved, and then continues again to increase until the next congestion event is detected. The initial implementation of TCP (i.e., *TCP Tahoe* [84, 86]), included the first three phases of congestion control, that is, the *slow-start*, *congestion avoidance*, and *fast recovery*. Whenever a loss was detected by way of receipt of three duplicate ACKs, congestion control was retransmitting the packet which was supposed to be lost, and then entered slow start phase. A subsequent update of the congestion control mechanism, known as *TCP Reno* [2], introduced the *fast recovery* phase, where the detection of loss via three duplicate ACKs does not require entering the slow start phase after *fast retransmit*, thus optimizing the

---

[1]A sender using SACK [35] must not send new data unless the incoming duplicate acknowledgment contains new SACK information.

performance by remaining in the *congestion avoidance* phase. The entire congestion control mechanism described earlier, is based on the *TCP Reno* IETF standard; however, several optimizations took place since then (e.g., *TCP NewReno* [40] which optimized the *fast recovery* phase, and the successor of them, *CUBIC* [89]). The following subsections describe in more detail how the legacy congestion control mechanism is operating in the context of MPTCP, as well as new algorithms and variations of the existing ones.

*2.1.10   MPTCP Implementation.* As mentioned earlier, in case of MPTCP, the application maintains a single socket for each individual MPTCP connection, regardless of the number of subflows associated with this connection. The life cycle of each individual subflow adheres to the legacy single-path TCP principles; consequently each subflow maintains the same mechanisms and operates similarly to a standard TCP connection. On the sending node, each subflow maintains its own subflow-level send window, as well as its own congestion and flow control mechanisms, which enable TCP-like communication with the receiving host at subflow-level. Thus, each subflow is individually responsible for the reliable and in-order delivery of the segments it carries. When a loss is detected, the subflow-level congestion control mechanism of the sending node undertakes the task of retransmitting the missing segment, while the receiving node ensures that the segments will be rearranged in the appropriate order at subflow-level. For this purpose, each individual subflow makes use of a dedicated out-of-order queue and a corresponding receive buffer.

However, as explained earlier as part of MPTCP flow control, even if it is ensured that each individual subflow delivers the segments carried over it reliably and in-order, it is not guaranteed that, when combining the ordered segments from multiple subflows, the result will be the original application bytestream. The MPTCP packet scheduler at the sending node could, for instance, split a data flow into segments traversing multiple subflows with heterogeneous path characteristics; then, at the receiving node, the in-order re-arrangement of the segments at subflow-level does not provide any guarantee for in-order arrangement *across* subflows. This task is assigned to the MPTCP connection-level out-of-order queue which ensures the reassembly of segments arriving from different subflows, reconstructing essentially the original bytestream in the MPTCP receive buffer, and delivering it to the receiving application. This is why MPTCP makes use of a new sequence space; the *Data Sequence Number (DSN)* and *Data Acknowledgment (Data-ACK)* at connection-level complement the subflow-level standard TCP's *sequence (SN)* and *acknowledgment numbers (ACK)*. It is a prerequisite that each segment be acknowledged both at subflow- and connection-level before being passed to the receiving application's socket; only then is the sending node's connection-level send window able to slide further to the right, allowing the transmission of more data.

If a segment is never acknowledged on a specific subflow, a subflow-level timeout occurs indicating severe congestion or complete failure of that specific subflow. In that case, the segment is copied into the retransmission queue which is read on priority by the scheduler, who injects the missing segment to another subflow as soon as an opening in the congestion window permits. All segments residing within the shared retransmission queue have higher priority over the segments within the MPTCP send buffer; only when all segments from the retransmission queue are scheduled, thus the retransmission queue becomes empty, does the scheduler pull segments from the shared send buffer.

A concise view of MPTCP protocol's end-to-end service delivery model, is depicted in Figure 3.

## 2.2   Packet Scheduling Algorithms

Research efforts have yielded a multitude of packet scheduling algorithms, with each one being tailored to specific application demands. Besides the target application, packet schedulers can also be classified according to the inner techniques they employ to achieve their goals; be it heuristics or machine learning (ML) algorithms. This work seeks
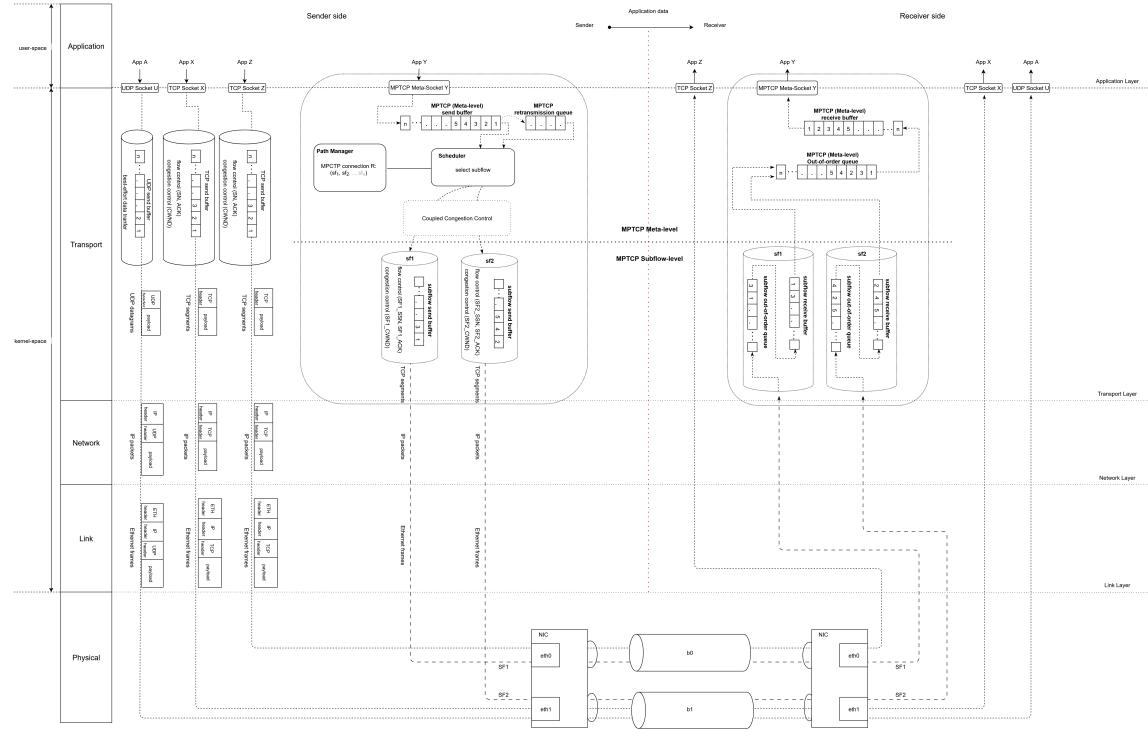
Fig. 3. Indicative MPTCP service delivery model

to present and assess the state-of-art *heuristics-based* MPTCP packet scheduling algorithms available in literature, up to the time of this publication. A more exhaustive list of multipath-compatible packet schedulers, along with some classification options, can be studied through other works available in literature, references to which are provided throughout the current manuscript. A brief overview of the MPTCP packet schedulers investigated through this work, is provided in Table 1.

*2.2.1 MPTCP Default (minRTT).* minRTT is the default MPTCP packet scheduler. It selects the subflow with the lowest RTT which has space in its congestion window (*cwnd*), to transmit the next segment(s). If the subflow experiencing the lowest RTT among the available subflows (i.e., among all subflows established and not marked as 'backup') does not have space in its cwnd, minRTT selects the subflow with the second lowest RTT that has sufficient cwnd. minRTT scheduler may optionally use some additional mechanisms to optimize its operation when the MPTCP connection is receive-window limited; that is, the receive window is blocked by one or more segments not yet acknowledged (e.g., missing segments or segments traversing a slow subflow), which are in turn causing HoL blocking at the receiving end. In such cases, *opportunistic retransmission and penalization* [79] technique showcases some additional benefit. *Opportunistic retransmission* mechanism is responsible for retransmitting a not-yet-acknowledged segment over a faster subflow, in anticipation for expedited delivery and unblocking of the receive window. In addition, the slow subflow is *penalized* by halving its congestion window and setting *ssthresh* to the new cwnd value. While those mechanisms provide some short-term advantage, their benefit is questionable in the long run; the reactive approach those mechanisms follow,

Table 1. MPTCP packet schedulers

| Scheduler | Full Name | Reference | Approach | Main concept |
|---|---|---|---|---|
| minRTT | default | v0.96 [6, 17, 18, 71, 79] | reactive | This is the default MPTCP packet scheduler. Assigns segments to the subflow experiencing the lowest RTT, as long as CWND permits. |
| BLEST | Blocking Estimation | v0.96 [6, 18, 29, 30] | proactive | BLEST tries to estimate and eliminate HoL blocking, by assessing whether it is beneficial to schedule segments on a currently available subflow, or skip it, waiting for a more advantageous one. |
| ECF | Earliest Completion First | v0.96 [6, 18, 63, 64] | proactive | ECF prioritizes scheduling over the faster subflows, anticipating to achieve lowest flow completion time. ECF may elect not to use a slower path, in an attempt to increase utilization of the faster ones. |
| RR | roundrobin | v0.96 [6, 18, 69, 71] | blind (ack-clocked) | Experimental scheduler. RR transmits traffic in a round-robin fashion. Tunable configuration following either real round-robin scheduling or ACK-clocked behavior. |
| LLHD | Low Latency and High Data Rate | v0.96 [66, 67] | reactive | LLHD bases its decisions on the combined measurement of real-time throughput, delay, and packet loss rate, rather than relying on a single parameter. LLHD adjusts scheduling reactively according to the perceived network conditions. |
| ReMP TCP | redundant | v0.96 [6, 18, 39] | blind | Sends each data segment on all available subflows. It guarantees lowest per-packet latency and increased reliability, at the expense of limited goodput and higher flow completion time. |

limits artificially the congestion window of the slow path, causing it to increase slowly until the next HoL blocking occurrence. This continuous penalization approach has a detrimental effect on the slow subflows, the bandwidth of which is significantly underutilized.

The algorithmic steps of *minRTT* packet scheduling algorithm are presented in detail in 1.

*2.2.2 BLocking ESTimation-based (BLEST).* BLEST adapts proactively its scheduling decisions by estimating whether injecting segments on a slow subflow would incur HoL blocking when the fast subflow becomes available again. It tries to limit, in essence, the amount of time a fast subflow remains underutilized due to receive window (and the mirror-effect of MPTCP send window) blocking. To achieve this, BLEST estimates the amount of data that would be sent on a fast subflow during the duration of a slow subflow (*RTTs*). If this calculated amount of data is larger than the usable MPTCP send window [28], then potential scheduling on the slow subflow would block the entire MPTCP send window. Instead, BLEST selects to skip the currently available slow subflow, waiting for the faster one to become available again (i.e., by delivering some segments and opening some space in its congestion window). Using this core idea, BLEST tackles the HoL blocking issue which is present during scheduling across heterogeneous networks or paths, at the expense of underutilizing the slower paths in case their use is deemed detrimental in the long run. BLEST uses minRTT algorithm 2.2.1 to identify the subflow experiencing the lowest RTT at any given moment before applying its sophisticated scheduling logic. A correction factor lambda ($\lambda$) adjusts the scale of X (see algorithmic steps in 2), controlling essentially how strict or relaxed will the decision on selecting the slower path be; an increase in the value of $\lambda$ induced by a HoL blocking observation, incurs more stringent requirements on the selection of the slow subflow $x_s$.

The algorithmic steps of *BLEST* packet scheduling algorithm are presented in detail in 2.

*2.2.3 Earliest Completion First (ECF).* ECF aims to reduce the completion time of a data flow traversing heterogeneous paths. It achieves this by ensuring that the completion of a flow is not delayed by transmission over a slower path. Instead, ECF prioritizes scheduling on the fastest (in terms of RTT) path as long as it is available for sending; that is, having available space in its congestion window. Unlike minRTT, ECF may elect not to send on a slower path and rather wait for a faster one to become available, if this is estimated to be beneficial for the total flow completion time. Thus, a slower path can only be used as long as it does not delay the completion of the entire flow. The estimation of completion times is calculated using not only path RTTs, but also an estimation of paths' bandwidth incorporated within the congestion window, and the size of the connection-level send buffer. ECF uses a margin $\delta$ which incorporates the standard deviation of each path's RTT (denoted by $\sigma$) to compensate for RTT and CWND variabilities, a hysteresis constant $\beta$ which limits frequent switchover between subflows, as well as parameter $k$ which holds the number of bytes remaining to be sent in order to complete the transfer of the data flow.

The algorithmic steps of *ECF* packet scheduling algorithm are presented in detail in 3.

*2.2.4 Round-Robin (RR).* The round-robin scheduler selects the subflows to schedule segments on, in a round-robin fashion [71]. While this technique might offer load-balancing and equal use of the available subflow resources in case of homogeneous paths, it is detrimental in case of highly heterogeneous paths. Segments scheduled in a round-robin fashion over paths with different network characteristics will end up out-of-order at the receiving node, causing HoL blocking, limiting the receive window, and delaying the whole data flow completion. Such conditions have a domino effect, increasing retransmissions, and stalling the entire connection. RR scheduler's default operation is *congestion-window limited*, that is, trying to fill in the entire congestion window of each subflow; as a result, RR becomes *ack-clocked* [53] in case of bulk data transfers, waiting, in essence, for ACKs to open up some space in the subflow's congestion window before the scheduler is able to inject new segments on that subflow again. However, the MPTCP development community has provided the option to tune RR's behavior through configurable parameters which control both the number of consecutive segments to be transmitted in each subflow before switching to another, as well as whether RR will leave some open space in the congestion window to retain its real round-robin scheduling scheme, avoiding the ack-clocked behavior [6].

The algorithmic steps of *RR* packet scheduling algorithm are presented in detail in 4.

*2.2.5 Low Latency and High Data Rate (LLHD).* Instead of relying on a single parameter to devise the next scheduling action, LLHD seeks to combine all the available network parameters (i.e., delay, path loss, and available bandwidth) and adapt quickly to the constantly changing network conditions. For this purpose, LLHD tries to optimize a utility function $\gamma$, which embodies the various network parameters as follows:

$$\gamma = GP_N + \beta \times \frac{1}{RTT_N} \tag{8}$$

where $\beta$ is a balancing factor, and $GP_N{}^2$ and $RTT_N$ are the normalized goodput and RTT, respectively, which are defined as follows:

$GP_N = \frac{GP}{GP_{max}}, \quad RTT_N = \frac{RTT}{RTT_{max}}$

where $GP$ and $RTT$ are the goodput and RTT of $sf_i$, respectively; similarly, $GP_{max}$ and $RTT_{max}$ are the maximum goodput and RTT values among all subflows, respectively. The algorithm considers available all those subflows with sufficient congestion window (cwnd). The subflow that maximizes the utility function $\gamma$ is the one to be selected for

---

[2] $Goodput = Throughput - Losses$

scheduling at any given time. Subflows that have been marked as "backup" are used only when no normal subflows are available. In addition, the algorithm examines whether a subflow is *unavailable* (e.g., no longer usable), or *temporarily unavailable* (e.g., no available space in its cwnd) for data transfer. The LLHD algorithm presented in 5, is based on the respective paper [67]. However, the set of steps is simplified; multiple conditions are examined within the actual code published in GitHub [66], controlling how the temporarily unavailable subflows, as well as any backup ones, are actually used upon each call of the algorithm.

The algorithmic steps of *LLHD* packet scheduling algorithm are presented in detail in 5.

*2.2.6 ReMP TCP (Redundant).* In multipath scheduling, redundancy is realized by sending duplicate segments on all available subflows. This technique is profoundly inefficient in terms of bandwidth utilization; it might however be applied to reduce per-packet latency and any induced retransmissions in case of interactive bandwidth-tolerant applications, which are though limited by stringent latency and jitter requirements. Using redundant packet scheduling, the fastest path is the one determining the effective end-to-end latency. In addition, redundancy contributes to reduced RTT, minimizing delays induced by excessive retransmissions; this can be achieved since segments arriving at the receiving node via the faster path will be acknowledged faster compared to how their duplicates (carried over a slower and even lossy path) would. According to its authors, *ReMP TCP* can efficiently handle queuing delays, latency variations (jitter), and heterogeneous path performance, exchanging, in essence, bandwidth for reduced latency in existing best-effort networks. If a subflow is not available at the time of redundant scheduling (i.e., in case cwnd does not permit, or in case of link unavailability or failure), a segment is injected only on the available subflow.

The algorithmic steps of *ReMP TCP* packet scheduling algorithm are presented in detail in 6.

## 2.3 Congestion Control Algorithms

*Uncoupled Congestion Control.* As with standard TCP, MPTCP also makes use of the congestion control mechanism to regulate the amount of outstanding data, in an attempt to detect and recover from network congestion. The standard TCP congestion control mechanism is applied individually to each MPTCP subflow, managing the rate at which each subflow can inject segments into the network. Thus, each MPTCP subflow maintains its own congestion window, operating essentially as an individual single-path TCP connection. While such an approach is functional, it leads to an unfair share of a bottleneck link's bandwidth. For instance, when two subflows belonging to the same MPTCP connection coexist alongside a standard single-path TCP connection in a bottleneck link, the MPTCP connection would occupy 2/3 of the bottleneck link's bandwidth. Applying such a congestion control mechanism for broader utilization, would provide an unfair advantage to MPTCP connections over standard TCP ones, allowing the latter to be allocated only a small fraction of the available bandwidth, and eventually leading them to resource starvation. A multipath congestion control mechanism with the aforementioned characteristics is called *uncoupled*, since the congestion window of each subflow scales individually, absent any MPTCP connection-level limit.

*Coupled Congestion control.* On the contrary, one of the main principles to be considered during the design of TCP congestion control algorithms, is to ensure *fairness* among the TCP connections which are competing for bandwidth [31]. The same principle is carried over to the MPTCP congestion control design [32, 36, 51, 78], ensuring that each MPTCP connection will be allocated the same amount of bandwidth as if it were a standard TCP connection; in the above example, the MPTCP connection would be allocated 1/2 of the bottleneck bandwidth, however many subflows comprising it. Such a congestion control algorithm, which imposes an upper bound on the aggregate amount of outstanding data at MPTCP connection-level, is called *coupled*.

Table 2 provides a concise description of the state-of-art congestion control algorithms used in MPTCP, inherited either directly from TCP (uncoupled), or adapted and new ones focusing on fairness (coupled).

*2.3.1 CUBIC.* The initial AIMD approach, as specified and applied in legacy *TCP Reno* and *NewReno* algorithms, was following a linear increase of the congestion window (cwnd) during the congestion avoidance phase, that is, an increase of the cwnd by one segment per RTT. While this approach is suitable for low BDP (bandwidth-delay product) links, there is a noticeable drawback when used in long-fat networks (i.e., with links experiencing high bandwidth and round-trip time delays). In such cases, the available bandwidth is underutilized since the linear increase of the congestion window takes too long to reach the link's saturation point, that is, the link's full capacity. As a result, the amount of time needed to reach the link's full capacity may be much higher than the one a flow may take to complete. For instance, using a 10Gbps/100ms round-trip delay link, the congestion window would reach the link's full capacity in $\sim 1.4$h.

The above computation is based on S. Floyd's et al. calculation of the round-trip times required until full congestion window is reached following a loss [34], that is:

$$\frac{\beta}{a}W + 1 \quad \text{round-trip times} \tag{9}$$

where $W$ is the congestion window value just before the loss event occurs, and $a$ and $\beta$ are AIMD's additive-increase and multiplicative-decrease factors, respectively. In the above example, using traditional TCP Reno and NewReno where $AIMD(a, \beta) = AIMD(1, 1/2)$, it would take approximately (given that: $1\,Byte \approx 10\,bits$, $MSS \approx 1000\,Bytes$):

$RTT * \frac{\beta}{a}W = 0.1 * \frac{1}{2}10^5 = 5000$ sec. for the congestion window to reach its full capacity. To cope with such long-distance and high-latency links, CUBIC replaces AIMD's additive increase alorithm with a cubic function, the window growth of which, is provided in equation 10:

$$W(t) = C(t - K)^3 + W_{max} \tag{10}$$

where $C$ is a CUBIC parameter determining the aggressiveness of the protocol against competing flows which may be using other congestion control schemes, $t$ is the elapsed time from the last window reduction, $W_{max}$ is the congestion window value when the loss occurred, and $K$ is the time period that the above function takes to increase $W$ to $W_{max}$ when there is no further loss event, and is calculated by the following equation:

$$K = \sqrt[3]{\frac{W_{max}\beta}{C}} \tag{11}$$

CUBIC heuristically sets constant $C = 0.4$, and the AIMD decrease factor $\beta = 0.2$, to adapt for fairness with other TCP flows, and attain a balance between protocol stability and convergence speed. CUBIC algorithm essentially computes the target value of the congestion window growth during the next RTT: $W(t + RTT)$, and then compares this target value with the growth of the congestion window that standard TCP would achieve within the same time interval. If the current congestion window is less than the estimated standard TCP's congestion window (calculated in terms of time $t$, where t is the elapsed time after the loss event), then CUBIC operates in *TCP mode* where the current congestion window is updated according to standard TCP's AIMD; this allows CUBIC to perform equally to TCP Reno in low BDP networks. If the current congestion window is greater than the one TCP would reach within the same time interval, then the protocol operates either in the *concave* or the *convex* region, depending on whether the current congestion window value is less or greater than $W_{max}$, respectively. Within both the concave and convex regions, CUBIC increments the congestion window (cwnd) by $\frac{W(t+RTT)-cwnd}{cwnd}$. The *concave* profile allows the algorithm to approach quickly the area

Table 2. MPTCP Congestion Control Algorithms (CCA)

| CCA | Full Name | Reference | Approach | Main concept |
|---|---|---|---|---|
| CUBIC | - | v0.96 [18, 41, 42, 89] | uncoupled (loss-based) | CUBIC replaces the linear additive increase phase of Reno-based AIMD CCAs, with a cubic function. This allows fast convergence toward the actual link bandwidth in high-BDP networks. |
| Coupled (LIA) | Linked Increases Algorithm | v0.96 [18, 70, 78, 88] | coupled (loss-based) | LIA is the default MPTCP congestion control algorithm. It has been the first coupled congestion control scheme to ensure MPTCP's incentive and fairness to single-path TCP flows. It modifies the additive increase part of Reno-based AIMD algorithm. |
| OLIA | Opportunistic Linked-Increases Algorithm | v0.96 [18, 57–59] | coupled (loss-based) | OLIA has been proposed as an optimization alternative to LIA, managing to eliminate the trade-off between responsiveness and optimal congestion balancing; thus, providing a pareto-optimal solution that combines both simultaneously. As with LIA, the proposed optimization is applied to the additive increase part of the AIMD algorithm. |
| BALIA | Balanced Linked Adaptation | v0.96 [18, 48, 74, 85] | coupled (loss-based) | BALIA reveals an inevitable tradeoff between friendliness, responsiveness, and window oscillation of congestion control schemes. Basing its design on an analytical framework, BALIA proves to reach a unique equilibrium point, managing to mitigate the trade-off and provide balanced performance. BALIA adjusts the entire AIMD algorithm controlling the Congestion Avoidance (CA) phase. |
| wVegas | weighted Vegas | v0.96 [8, 9, 18, 90] | coupled (delay-based) | wVegas is a coupled delay-based alternative of standard TCP-Vegas for MPTCP. Contrary to loss-based CCAs that associate packet loss with congestion, wVegas interprets queuing delay as congestion signal. Using weights to equally distribute the congestion cost among subflows, wVegas ensures that each sublow allocates a fair share of the bottleneck bandwidth. Maintaining a moderate buffer occupancy in link queues, wVegas accounts for congestion proactively, thus preventing excessive packet loss. |
| BBR | Bottleneck Bandwidth and Round-trip propagation time | v0.96 [10–12, 16, 18, 54] | uncoupled (congestion-based) | BBR-v1 [11] is neither a loss- nor a purely delay-based scheme. Although inspired by the delay-based approach of TCP-Vegas, BBR follows a unique method to identify the optimal operating point, by accurately estimating the bottleneck bandwidth and base RTT. Based on these estimations and appropriate traffic pacing adjustments, BBR achieves maximum bandwidth utilization at the lowest possible round-trip delay, eliminating intermediate queue formation. To derive these estimations, BBR implements a state transition model, cycling over four states (Startup, Drain, Probe_BW, and Probe_RTT). BBR's congestion-proactive approach leads to sustained high bandwidth utilization, low RTT, and infrequent loss. |
| C-MPBBR | Coupled Multipath BBR | v0.96 [55, 68] | coupled (congestion-based) | C-MPBBR is a coupled alternative of BBR for MPTCP. It exploits BBR's existing BtlBw and DelRt parameters to ensure multipath benefit, as well as fairness among single- and multi-path flows traversing a common bottleneck. C-MPBBR achieves the first goal by closing any subflows not contributing sufficiently to the multipath connection. Fairness is established by leveraging BtlBw estimations to identify subflows on the same bottleneck, and dividing the measured bandwidth equally among them. |

close to link's saturation point, but then slow down to remain long-enough in this area which is expected to provide the highest bandwidth utilization. In case network conditions have improved since the last congestion event (i.e., the reception of a duplicate ACK), the protocol is switching to the *convex* region (also known as *maximum probing phase*) where it seeks to identify the new bandwidth saturation point; it does this following the convex function's exponential nature, in which the algorithm prompts for the new $W_{max}$ initially slowly, and then more aggressively.

The advantage of CUBIC compared to standard TCP is that, while maintaining the standard TCP principles, it manages to converge faster in case of high BDP links. It is also unique in the way it approaches the saturation point, by combining the concave and convex profiles encompassed within a single cubic function. This way, it approaches quickly the last epoch's $W_{max}$, then prompts slowly for the new saturation point close to the plateau created around $W_{max}$, and in case the available bandwidth has increased since the previous epoch, it increases the congestion window's growth rate more aggressively to find the new maximum value. Another key property of CUBIC is the fact that the congestion window increase rate is independent of RTTs outside the TCP-mode; this characteristic allows flows with different RTTs to maintain similar congestion window sizes, thus further contributing to the protocol's fairness.

Since CUBIC is based on the AIMD congestion control scheme, where the congestion window decreases[3] rapidly upon an identified loss and then grows linearly, it is considered a *loss-based* congestion control algorithm. From MPTCP perspective, since CUBIC does not incorporate any multipath-specific logic to account for fairness between MPTCP and standard TCP connections, it is considered *uncoupled*. A condensed version of CUBIC pseudo-code is presented in the Appendix, based on the original source [42]. Some more recent updates can be found within RFC 9438 [89].

The algorithmic steps of *CUBIC* congestion control algorithm are presented in detail in 7.

### 2.3.2 *Coupled.*

Coupled (also called *Linked Increases Algorithm (LIA)*) [78, 88], has been the first coupled congestion control algorithm for Multipath TCP that efficiently tackles the fairness issues arising when single- and multi-path flows concurrently share the same bottleneck link resources. Coupled congestion control has been therefore introduced to satisfy the fairness principles that the respective standard TCP congestion control schemes lack.

The bottleneck fairness properties that any multipath congestion control scheme should employ, are succinctly summarized in the following three goals [78]:

- Goal #1 (Improve Throughput): A multipath flow should at least achieve the performance of a single path flow on the best of the paths available to it.
- Goal #2 (Do no harm): A multipath flow should not allocate more capacity from the resources shared when traversing multiple paths, than if it were a single flow using only one of these paths. This guarantees fairness against other competing flows.
- Goal #3 (Balance congestion): A multipath flow is expected to remove traffic from its most congested paths, subject to meeting the first two goals.

On the basis of these goals, three main approaches have been studied in literature in an attempt to identify the most appropriate multipath congestion control coupling scheme [44]. The first approach suggests the use of a common congestion window, shared between the subflows of a multipath connection. While this approach sounds reasonable, it may lead to suboptimal performance of the multipath connection, since individual subflows traversing distinct paths usually experience different RTT characteristics, which in turn lead to varying ACK receipts and, consequently, to

---

[3]To avoid any confusion between the multiplicative decrease factor as defined within [42] ($\beta = 0.2$) and the one specified within [89] ($\beta = 0.7$), we have to clarify that both are using the same decrease factor. The RFC updates *ssthresh* by exploiting the decrease factor directly, i.e. by: $ssthresh = cwnd * \beta$, while the paper calculates it using: $ssthresh = (1 - \beta) * cwnd$.

differing congestion window increase and decrease intervals. As a result, there is no perfect synchronization among subflows regarding the congestion window increase and decrease events. Furthermore, a packet loss event on a particular subflow would decrease the shared congestion window, degrading the performance of the entire multipath connection. The second approach includes the detection of shared bottleneck links, and ensuring that no more than one subflow will traverse such a link. While that approach was not that convincing a couple of years ago [44, 92], recent works provide the means for accurate detection of shared bottlenecks [10], enabling also efficient coupling approaches [68]. The third approach considers individual congestion windows for each subflow of a multipath connection, where each subflow adjusts its own congestion window according to a weight, such that the entire multipath flow (comprising multiple subflows) has the same aggressiveness as a standard TCP flow [44]. Building upon this approach, the work pursued in [44, 88] suggests to modify the legacy additive increase algorithm of the congestion avoidance phase, as follows:

- upon receipt of an ACK, increase congestion window by: $\frac{\alpha}{w_r}$
- upon loss detection, decrease the congestion window by: $\frac{w_r}{2}$

where $w_r$ is the congestion window of subflow $r$. According to [44], the relation between the increase parameter $\alpha$ and each subflow's individual weight factor $D$ is denoted by $\alpha = D^2$ (a). A weight factor of $D = \frac{1}{n}$ (b), where $n$ is the number of subflows belonging to a multipath connection, would couple the throughput of a multipath connection to the equivalent of a standard TCP connection within an RTT interval. Consequently, according to (a) and (b), the increase parameter should be: $\alpha = \frac{1}{\sqrt{n}}$. While this approach ensures fairness among single- and multi-path flows traversing a shared bottleneck link where flows experience the same RTTs, it may lead to suboptimal resource utilization when the flows traverse multi-hop bottleneck paths encompassing varying RTTs and packet loss rates [88].

This is where *Coupled* congestion control algorithm kicks in, managing to equalize the multipath flow's aggregate bandwidth with the one a regular TCP flow would get on the best path available to the multipath flow. Coupled algorithm enables the multipath flow to first estimate the target rate of a regular TCP flow based on currently measured path conditions (i.e., RTT, packet loss rate). Then, it calculates the increase parameter $\alpha$ which controls the overall aggressiveness of the protocol and assists in reaching the desirable target rate. Similarly to EWTCP [88], Coupled modifies the additive increase algorithm of the congestion avoidance phase, while maintaining the legacy TCP's multiplicative decrease, fast retransmit, and fast recovery schemes [2, 78].

More specifically[4],

- upon ACK receipt on $subflow_r$, increase the congestion window $w_r$ by: $\min(\frac{\alpha}{w_{total}}, \frac{1}{w_r})$
- upon loss detection on $subflow_r$, decrease the congestion window $w_r$ by: $\frac{w_r}{2}$

Noteworthy here, is the fact that Coupled algorithm's design encompasses an adaptive behavior based on the level of statistical multiplexing. In cases of high statistical multiplexing, where the multipath flow coexists alongside single-path flows, the multipath flow does not influence the link loss rates; in such scenarios, MPTCP will achieve the same throughput as standard TCP on the best path. However, when low statistical multiplexing is present, in which case the multipath flow influences the loss rate of the path, MPTCP's throughput is guaranteed to be higher than that of a single-path TCP flow on the best path; in that case, the achieved multipath throughput is expected to be equal to the sum of the idle paths' individual throughput, thus, fully aggregating the available bandwidth.

Contrary to EWTCP [88] where the increase parameter $\alpha$ was inversely proportional to the square root of the number of subflows, Coupled bases its decisions on how TCP would perform on the best path available at each RTT interval;

---

[4]Details on the specific equations which determine the congestion window increase and decrease steps, are provided within the algorithm in 8.

since path conditions in real networks are subject to constant and dynamic changes, Coupled needs to compute $\alpha$[5] once per RTT or upon packet loss detection[6][78].

Since Coupled CCA couples only the congestion window increases, it satisfies Goals #1 & #2. To achieve Goal #3, which essentially implies perfect resource pooling [87], would require to couple also the multiplicative decrease phase, while ensuring that no traffic would be scheduled on links with higher loss rates. According to LIA's authors, given the Internet's dynamic path conditions, such a restriction would lead the algorithm to *flappiness* between paths and insufficient probing of lossy paths, thus failing to utilize them quickly once they recover. The notion behind LIA's operation is to allocate congestion windows to the subflows, such that $p_i * cwnd_i$ is constant, for all $i$ ($p_i$ denotes the packet loss rate experienced by $sf_i$). This ensures that equal congestion windows will be allocated among subflows in case they experience equal loss rates, while higher congestion window values will be progressively allocated to less-lossy subflows when their packet loss rates differ.

The algorithmic steps of *LIA* congestion control algorithm are presented in detail in 8.

### 2.3.3 Opportunistic Linked-Increases Algorithm (OLIA).

LIA's shortcoming in fulfilling the third principle of multipath congestion control, has gathered attention. Some more complex scenarios examining MPTCP connections' fairness to regular TCP flows under LIA congestion control, identified a couple of issues: i) the fact that upgrading TCP users to MPTCP has a negative impact to the remaining regular TCP flows without any added benefit to the MPTCP ones, and ii) that MPTCP is over-aggressive to regular TCP flows [59]. Scientific research on this field, revealed that it is feasible to attain all three multipath congestion control goals simultaneously. As a result, *OLIA* congestion control algorithm has been proposed as an alternative to LIA. Contrary to LIA, OLIA's design eliminates the trade-off between responsiveness and optimal congestion balancing, managing to provide a pareto-optimal solution which achieves both simultaneously.

Similarly to LIA, OLIA couples only the additive increase algorithm of the congestion avoidance phase; in case of loss, the multiplicative decrease, as well as the fast retransmit and fast recovery, retain the legacy TCP NewReno approach.

OLIA maintains a couple of variables, such as the set of paths over which subflows of the MPTCP connection have been established (*all_paths*), a subset of "all_paths" containing all presumably best paths (*best_paths*), a subset of "all_paths" including paths with the largest congestion windows (*max_w_paths*), and a subset of "all_paths" which also belong to "best_paths" but are not among those exhibiting largest congestion windows (*collected_paths*). Exploiting these variables, the increase part of OLIA is defined as follows:

- for each ACK received on path r, increase the congestion window ($w_r$) by:

$$\frac{w_r/rtt_r^2}{\left(\sum_{p \in R_u} w_p/rtt_p\right)^2} + \frac{a_r}{w_r} \tag{12}$$

  where $p$ denotes a path on the set of "all_paths" $R_u$, and $w_p$ and $rtt_p$ denote the congestion window and the round-trip time of path $r$, respectively. Equation 12 comprises two terms: the first term ensures Pareto optimality in resource pooling (satisfying MPTCP Goal #3: "Balance congestion"), while the second one, encompassing $\alpha_r$, guarantees responsiveness and non-flappiness (satisfying MPTCP Goal #1: "Improve throughput", and Goal #2: "Do no harm") [59]

---

[5]$alpha$ calculation formula: The calculation of $\alpha$ within algorithm's step #2 in 8, is derived by equalizing the rate of the multipath flow with the rate of a TCP running on the best path, and solving for $\alpha$.
[6]The pseudo-code steps showcase the congestion window increment and $\alpha$ computation per ACK; in actual implementations, however, such computations are performed once per RTT or upon packet loss detection.

$a_r$ is calculated as follows:

- if $r$ is within the *collected_paths*, then: $a_r = \frac{1/num\_of\_paths}{|collected\_paths|}$
- if $r$ is within the *max_w_paths* and if *collected_paths* is not empty, then: $a_r = -\frac{1/num\_of\_paths}{|max\_w\_paths|}$
- otherwise, $a_r = 0$.

- for each loss on path r, decrease the congestion window ($w_r$) by: $\frac{w_r}{2}$.

OLIA adapts the congestion window increases as a function of the underlying paths' RTTs, thus efficiently compensates for different RTTs. Analytical results and testbed experiments prove that OLIA provides fairness among coexisting TCP and MPTCP flows, as well as optimal congestion balancing, while remaining as responsive and non-flappy as LIA.

The algorithmic steps of *OLIA* congestion control algorithm are presented in detail in 9.

### 2.3.4 Balanced Linked Adaptation (BALIA).

A subsequent work has identified that OLIA also suffers a shortcoming, that is, its limited responsiveness to dynamic network conditions when the used paths experience similar round-trip times (RTTs). To confront OLIA's weakness, meticulous research conducted in the context of [74] devised BALIA which, similarly to its predecessors LIA and OLIA, is a window-based coupled congestion control algorithm designed for MPTCP. This work identified an inevitable tradeoff between friendliness, responsiveness, and window oscillation of congestion control schemes, and proposed an analytical framework that enables the aforementioned properties' assessment for any congestion control algorithm, already from the design phase.

Leveraging the designed analytical framework to prove BALIA's unique equilibrium point and asymptotical stability, and confirming it through experimental evaluation, BALIA's authors managed to successfully mitigate the tradeoff and provide balanced performance. As with LIA and OLIA, BALIA modifies the additive increase (AI) algorithm of the congestion avoidance phase, while maintaining the TCP NewReno fast retransmit & recovery algorithms. In case of loss, BALIA adjusts the TCP NewReno multiplicative decrease (MD) algorithm, multiplying it by a factor in the range of [1, 1.5]. In case of single path use, BALIA reduces to the standard TCP NewReno increment and decrement algorithms.

The algorithmic steps of *BALIA* congestion control algorithm are presented in detail in 10.

### 2.3.5 Weighted Vegas (wVegas).

An alternative approach to the aforementioned loss-based congestion control schemes (i.e., LIA, OLIA, BALIA), is *wVegas* which bases its decisions on queuing delay measurements. wVegas builds upon the respective standard TCP congestion control algorithm TCP-Vegas, and is adapted to serve the fairness requirements of multipath environments. Contrary to other coupled multipath congestion control schemes (e.g., LIA, OLIA, and BALIA) that estimate network congestion through packet losses, wVegas interprets queuing delay as congestion signal, thus being more sensitive to network condition changes in an attempt to achieve fine-grained load balancing. wVegas focuses on satisfying the "Congestion Equality Principle" [9], where all flows attempt to allocate a fair share of the available bottleneck bandwidth via an equal share of the congestion cost. For this purpose, wVegas assigns a weight to each subflow belonging to a multipath flow, and adaptively adjusts this weight according to the Congestion Equality Principle. The weight quantifies, in essence, the aggressiveness of each subflow on acquiring part of the available bandwidth. Subflows established on less congested paths are expected to be assigned a larger weight value, allowing them to compete for bandwidth more aggressively; this competition is expected to increase the underlying path's congestion, which will in turn lead to decrease of the subflow's weight value on that path. This cycle is executed for all subflows traversing the various paths until congestion equilibrium is reached; at that point network resources will be fairly shared by all flows.

Diving deeper into the design principles of the algorithm, wVegas is essentially assigning a fixed parameter $\alpha = 10$ to each one of the multipath flows. This parameter determines the total number of bytes backlogged in the network for all subflows belonging to a MPTCP flow [9, 90]. Each subflow (e.g., assuming a subflow on path $r$) belonging to a broader multipath flow, is then assigned a portion of this value (e.g., $\alpha_r$) according to the weight corresponding to this subflow $w_r$. The weight is calculated as the portion of a subflow's rate to the overall rate of the entire multipath flow, as depicted in 14.

The *rate* and *weight* of a subflow $r$ is respectively provided by the below equations:

$$rate_r = cwnd_r / rtt_r \tag{13}$$

$$weight_r = rate_r / \sum_i (rate_i) \tag{14}$$

where $\sum_i (rate\_i)$ denotes the aggregate rate achieved by all subflows of the multipath flow.

The portion of $\alpha$[7], which then corresponds to a subflow on path $r$, $\alpha_r$, is calculated as follows:

$$\alpha_r = weight_r * \alpha \tag{15}$$

Within wVegas, $\alpha_r$ serves as a threshold which controls whether the equilibrium rate should be updated, and whether the congestion window of the respective subflow established on path $r$ will increase or decrease at the end of each RTT interval. If the currently measured rate deviates from the expected equilibrium rate beyond $\alpha_r$, then the congestion window is decreased at the end of the RTT; otherwise it is increased. The aforementioned rate difference is calculated within the wVegas' weight adjustment algorithm as follows:

$$diff = \left( \frac{cwnd}{baseRTT} - \frac{cwnd}{rtt} \right) \cdot baseRTT \tag{16}$$

Consequently, wVegas adjusts parameter $\alpha$ dynamically, thereby influencing the transmission rate of the corresponding subflow and managing to equalize congestion on the path. Normalized $\alpha$, that is $\frac{\alpha_r}{\alpha}$, is essentially the weight factor $w_r$ applied on the subflow established on path $r$. The packet queuing delay measurement employed by wVegas is meant to provide a more accurate estimation of the underlying path's congestion compared to loss-based techniques. Since wVegas maintains a moderate buffer occupancy in link queues, it accounts for congestion proactively and rarely leads to packet losses. On the contrary, loss-based congestion control schemes try to fill in intermediate queues before a loss is detected and backoff procedure begins. However, wVegas experiences some shortcomings, such as its reliance on accurate RTT measurements, its less aggressive behavior when competing with loss-based algorithms, as well as its limited efficiency on high BDP paths.

The algorithmic steps of *wVegas* congestion control algorithm are presented in detail in 11.

*2.3.6 Bottleneck Bandwidth and Round-trip propagation time (BBR).* Traditional loss-based congestion control algorithms cycle through the process of filling in the bottleneck link and any intermediate queues, and then responding to the detected bufferbloat-induced loss. Deviating from such loss-based schemes, *BBR* follows a different approach, inspired by the delay-based control logic proposed within TCP-Vegas [5]. Contrary to the operation of loss-based schemes, BBR attempts to identify a bottleneck link's optimal operating point. Such a point entails maximum utilization of the underlying link's capacity at the lowest round-trip delay. At each moment, a sender operating at this optimal

---

[7]Parameter $\alpha$ is associated with a MPTCP flow, while $a_r$ is associated with a specific subflow $r$ of that MPTCP connection. Parameter $\alpha$ is referred to as *total_alpha* within the algorithmic steps in 11. Similarly, the aggregate rate of all MPTCP subflows, denoted as $\sum_i (rate_i)$ in 14, is referred to as *total_rate* into the algorithmic steps.

point, transmits at the bottleneck link rate, while simultaneously ensuring that no intermediate queues are formed. This prevents round-trip delay increases and minimizes the packet loss probability.

For this purpose, BBR makes use of two parameters: *BtlBw* (bottleneck bandwidth) and *RTprop* (round-trip propagation time); the first tracking bottleneck link's capacity, and the latter the round-trip delay. Rtprop reflects the base RTT and is used to measure any subsequent round-trip delay updates. Given a fixed data path, RTprop remains constant up to the point where intermediate queues begin to be formed; from that point onward, an increase in RTT indicates that the pipe is already full and the sending rate exceeds bottleneck bandwidth. An increase in RTprop while BtlBw remains constant, allows BBR to infer that intermediate queues have started to be filled in. Bottleneck bandwidth (BtlBw) is used to probe for any updates in the underlying bottleneck capacity by comparing it with the actual delivery rate (DelRt). BBR attempts to estimate the bottleneck link's optimal operating point by trying to estimate these two parameters. The constraints reflected through these parameters are intrinsically contradicting, meaning that the operating point where one can be measured, hinders the measurement of the other. Measuring RTprop, for instance, requires that inflight packets remain below BDP, so that the link is underutilized. On the other hand, bottleneck bandwidth estimations request for probing above the known BDP, which could lie in the region where queues start to be formed, in case bottleneck capacity has not changed.

In order for a link to be fully utilized at each moment, BDP amount of data should be inflight, where $BDP = BtlBw * Rtprop$. BtlBw controls the sending rate and ensures full bandwidth utilization, while Rtprop guarantees an always full pipe by managing the volume of inflight data; an amount of data less than the optimal would cause link starvation, while an excess would overfill the pipe and lead to developing queues. To accurately measure BtlBw and RTprop, BBR implements a state transition model, cycling over different states, such as the *Startup*, *Drain*, *Probe_BW*, and *Probe_RTT*. During Startup phase, a flow's congestion window increases exponentially, similarly to that of loss-based congestion control schemes; however the congestion window increase with BBR is smoother and the algorithm does not back off in case of packet loss or delay increase, maintaining its robustness toward discovering the available bandwidth. The Startup phase is over once the bottleneck bandwidth is considered to be successfully detected. Then, BBR enters the Drain phase, the purpose of which is to drain any queues formed during Startup. The Drain phase lasts one RTT interval, managing to clean up any filled queues rapidly. Once BBR estimates that the queue is empty and the link is sufficiently filled with BDP data inflight, it leaves Drain and enters Probe_BW phase. BBR consumes the majority of its time in this phase, keeping a BDP of data inflight, paced at a rate equal to the BtlBw estimate. The *pacing_gain* parameter used in this phase controls the way the algorithm probes for updates in BtlBw estimation. More specifically, Probe_Bw phase consists of eight-cycle rounds, each lasting an RTprop interval. A full Probe_BW round involves cycling through pacing_gain values in the following order: {5/4, 3/4, 1, 1, 1, 1, 1, 1} [10]. According to this pacing_gain sequence, BBR spends one round exploring whether bottleneck bandwidth has increased, thus probing at a rate equal to 5/4 of the current estimate. Then, pacing_gain decreases equally to the initial increase portion (3/4), draining any queues formed during the previous round. BBR spends the next six rounds "cruising" at BtlBw rate, thus achieving the highest link utilization with minimum delay.

If probing for higher bandwidth leads to higher delivery rate, BBR updates its bottleneck bandwidth estimation and adapts traffic transmission rate at the new delivery rate. If probing for higher bandwidth is not accompanied by delivery rate increase, it means that the current transmission rate already matches the bottleneck link's capacity, and the excess data transmitted during BtlBw probing leads to queue formation, which is also perceptible through a subsequent increase in RTprop.

In case RTprop estimate has not been updated by a lower RTT measurement within 10 seconds of operation into the Probe_BW phase, BBR enters Probe_RTT phase. The intention of this phase is to obtain an accurate measurement of the base RTT by reducing the congestion window (cwnd) to a very small value (i.e., four packets). After remaining in Probe_RTT phase for at least 200ms and one RTT [10], BBR enters either Startup or Probe_BW phase, depending of the BtlBw estimate; if BtlBw indicates a full pipe, BBR transitions to Probe_BW phase; otherwise to Startup.

Another important parameter used in BBR[8]is the *cwnd_gain* which bounds the amount of data inflight in each state; cwnd_gain is different from the pacing_gain, in that the first allows for an upper bound of inflight data, while the latter controls the sending rate and consequently influences the amount of outstanding data occupying the link.

What distinguishes BBR from AIMD algorithms, is the fact that it is proactive to loss, operating at the point where intermediate queues start to be formed, whereas AIMD schemes operate at the saturation point where queues are already full, merely reacting to detected loss. BBR's subtle traffic pacing at a rate equal to bottleneck link's bandwidth, differs from traditional AIMD algorithms whose sending scheme follows a bursty pattern within each RTT interval [47].

In the context of multipath communications, BBR is deemed an uncoupled congestion control scheme. In the following subsection, a coupled alternative, based on BBR, is presented (C-MPBBR).

The algorithmic steps[9]of *BBR* congestion control algorithm are presented in detail in 12.

*2.3.7  Coupled Multipath BBR (C-MPBBR).* C-MPBBR arose from an attempt to create a coupled congestion control algorithm for MPTCP, on the basis of BBR. The main goals C-MPBBR aims at fulfilling, are the following:

   i. Goal #1: to incentivize the use of MPTCP over standrard single-path TCP connections, which implies that MPTCP performs at least as well as a single-path TCP (SPTCP) flow on the best path, and
   ii. Goal #2: to reinforce fairness among flows consuming the resources of a shared bottleneck.

To achieve Goal #1, C-MPBBR continuously measures the bottleneck bandwidth (BtlBw) and the actual delivery rate (DelRt) using BBR's existing parameters. If the algorithm identifies that, during five consecutive ProbeBw rounds, the aggregate throughput of the MPTCP subflows is below that of a single-path flow on the best path, it closes the subflow experiencing the lowest performance. The algorithm performs the same operation for as long as the same condition holds true, up until the point where only a single subflow remains available; at that point, an MPTCP connection comprising one subflow is equivalent to a SPTCP flow, acquiring and utilizing the same amount of resources.

C-MPBBR adheres to Goal #2 by ensuring that the aggregate amount of bandwidth allocated to the subflows of a MPTCP connection equals that of a SPTCP flow, when sharing the same bottleneck link. C-MPBBR satisfies fairness among single- and multi-path flows, by exploiting the BtlBw parameter value. The BtlBw estimations, which are measured individually by each subflow, allow the algorithm to infer whether the subflows of a MPTCP connection traverse the same bottleneck link. In that case, the BtlBw value measured by an individual subflow is divided by the number of subflows sharing the same bottleneck link, so that multipath flows are allocated their fair share when competing against single-path flows for bottleneck resources.

C-MPBBR is a worth-researching coupled congestion control scheme, as it manages to combine the state-of-art technology inherited from BBR, along with MPTCP's *aggregation benefit* [56] and *fairness* concepts [27, 31]; initial tests [68] showcased that it fulfills both goals when compared to coupled loss-based congestion control schemes (i.e., LIA, OLIA, BALIA), or other BBR-based variants. For this reason, we selected to include C-MPBBR within our own

---

[8]It is important to mention here that throughout our work, any reference to BBR is referred to *BBR v1* [10, 11], which is also available in the older MPTCPv0 linux kernels. However, BBR is still under development at the time of this publication. IETF details in regard to BBR are referred to BBR v2 [12].
[9]An abstract version of BBR algorithm's pseudo-code is provided in 12, as originally published for BBR_v1 within [11]. BBR actually spans more than a thousand lines of C code, details of which can be found within [11, 12, 14, 54].

experimentation campaign, and verify its behavior both compared to other congestion control schemes, as well as in combination with various MPTCP packet scheduling algorithms.

The pseudo-code determining *C-MPBBR's* operation toward fulfilling its incentive and fairness goals, is presented in 13.

## 3 Experimentation Methodology

This section provides details on the methodology followed which laid the groundwork for conducting the experiments. It initially presents the multipath test-environment's topology, the tools used, and the configuration applied; it then continues with the experimentation campaign's design and the set of experiments executed. Finally, it delineates the main performance metrics against which the packet scheduling and congestion control algorithms have been assessed.
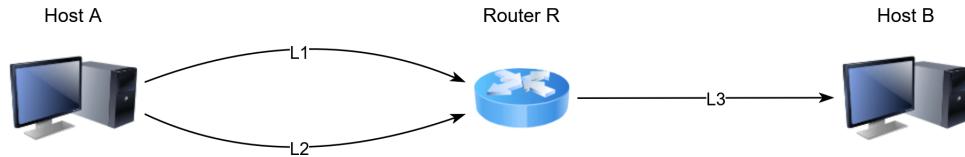
### 3.1 Multipath Environment Setup & Configuration



Fig. 4. Multipath test environment topology

*Mininet topology.* Figure 4 depicts our Mininet-based multipath topology comprising two multipath-capable hosts (A, B) which are connected via an intermediate router (R). More specifically, source host (i.e., Host A) is connected to R over two disjoint paths (L1, L2), and R is then connected to the destination host (i.e., Host B) over path L3. The MPTCP connection established between Hosts A and B, comprises two subflows (SF1, SF2), that is, one subflow is established over paths L1, L3 (i.e., SF1), and another subflow over paths L2, L3 (i.e., SF2). The two subflows are partially disjoint, since both are eventually traversing the same path L3 towards destination host (Host B).

The aforementioned topology has been selected based on the *hybrid-access* network as depicted within [3] (Fig.2: *Hybrid Access Network*), in an attempt to have a simple and explainable configuration that guarantees reliable performance results for the algorithms under study. In our hybrid-access topology, subflows SF1 (over paths L1, L3) and SF2 (over paths L2, L3) are established between the two interfaces available in Host A (e.g., eth0, eth1) and the single interface in Host B (e.g., eth0).

*HW/OS details.* The above-mentioned multipath topology is emulated using Mininet v2.3.0 [21, 25, 62] and Python v3.10 [77]. The Mininet environment is instantiated on top of Linux 22.04.5 LTS (Jammy Jellyfish) OS [7][10], which is in turn running on a Dell OptiPlex-7050 server, equipped with an Intel Core i5-6500 CPU @ 3.20 GHz, 16 GiB RAM, and 238 GiB SSD. To be able to execute Mininet with MPTCP capabilities, we loaded the out-of-tree MPTCP linux kernel v5.4.230 [38] (MPTCP v0.96 [18]) which we have further customized in order to be able to collect kernel parameter values, such as the MPTCP out-of-order queue [19, 20]. Another reason that forced us patch the legacy MPTCP kernel, was the fact that we needed to port all packet scheduling and congestion control algorithms to the same MPTCP version

---

[10]Noteworthy here is the fact that the deployment and execution of Mininet in a VM produced weird results, non-deterministic under identical test iterations, which forced us to execute Mininet directly on the OS.

(v0.96) for compatibility and fair comparison purposes; to this end, some modifications were also required to the legacy kernel code basis.

*MPTCP version.* The MPTCPv0 linux kernel version, which we selected to use for our experiments, has been recently frozen. MPTCPv1 efforts are ongoing to develop MPTCP in compliance with the latest IETF RFC [36]. The reason behind selecting the older MPCTP linux kernel version was the fact that the out-of-tree version (MPCTPv0 [22]) was mature enough and already included most of the MPTCP packet scheduling and congestion control schemes as loadable kernel modules. The new upstream version (MPTCPv1 [23]) facilitates algorithms' development in user- rather than kernel-space (e.g., via eBPF); thus, MPTCPv1 lacked the needed algorithms, offering only the default minRTT packet scheduler.

*MPTCP configuration.* Delving into TCP configuration details, MPTCP can be enabled via the `mptcp_enabled` parameter, and the desired packet scheduling and congestion control algorithm can be adjusted using the `mptcp_scheduler` and `tcp_congestion_control` linux parameters, respectively [6]. In our setup, the MPTCP default path manager, *fullmesh* has been used throughout the entire experimentation campaign (configured via `mptcp_path_manager` parameter). In addition, we have disabled `mptcp_checksum` since our Mininet experimentation platform operates in a controlled environment, absent any middleboxes which could potentially alter the payload [36]. Furthermore, `tcp_no_metrics_save` has been enabled, flushing metrics from the route cache once a connection is closed [43, 75]. Except these, all the remaining TCP parameters maintained their default values, allowing a performance evaluation as closest to a normal setup as possible; thus, TCP autotuning, and selective acknowledgement (SACK) features have been kept enabled. Noteworthy here is the fact that autotuning has been experimentally associated with MPTCP performance degradation and lower aggregation benefit, as well as with slower congestion window increase within slow-start phase, when compared to the equivalent manual setup to the same maximum buffer limit [72]; however, we have selected to keep it enabled in our environment to emulate the default setup's use case. Another important condition to be considered in any such multipath performance evaluation scenario, either emulated or real-world, is to guarantee that the end-hosts are not receive-window limited. This means that the send and receive buffers should be allocated enough memory so that flow control is not the limiting factor in the experiment, but rather the congestion control algorithm via the congestion-window. Thus, it has been suggested that, for MPTCP, the maximum send and receive buffer size be set to a value of $2 * BDP$, that is, $2 * \sum_{i=1}^{n} bw_i * RTT_{max}$, where $bw_i$ is the bandwidth of subflow $i$, and $RTT_{max}$ the highest RTT value measured across all subflows [1, 72, 79]. The logic behind this formula is that one BDP of memory is required to keep sending while waiting for a packet sent on the slowest path to be delivered; to keep sending when a path is used for fast retransmission, another BDP of buffer size is required [79]. In our system, where autotuning has been enabled, the maximum send and receive buffer size (wmem and rmem, respectively) have been set by default to 16MiB which is considered to be sufficient for our experiments and in accordance with the proposed value of *at least* $2 * BDP$ [1]. Details on the exact bandwidth, delay, and packet loss rate value ranges of each scenario is provided later in this section, but for a quick calculation, our experiments consider two subflows with a maximum capacity of 100 Mbps each, an $RTT_{max}$ of 50ms, and a maximum packet loss rate of 2%; these values require a buffer size of at least 1.25 MiB, which is covered by the existing autotuning upper limit.

*Loadable modules.* The packet scheduling and congestion control algorithms can be loaded in the linux kernel at runtime using `insmod` or `modprobe` commands, without the need to modify, rebuild, and install a new kernel. As a result, some work was initially required to port all algorithms to the common running-kernel basis (i.e., MPTCP v0.96

[18]). Consequently, we had to modify the legacy code for LLHD packet scheduler [66], and C-MPBBR congestion control algorithm [55].

*Additional tools.* Concerning the set of tools used to conduct the experiments, we have exploited `iperf` v2.1.5 to generate traffic between the two Mininet hosts [49, 50], `matplotlib` v3.10.0 for visualization of the collected results [45], as well as the `tcp_probe` linux feature which allowed us to record the MPTCP connection state and retrieve kernel metrics during test execution [15, 37].

### 3.2 Traffic methodology

Using iperf tool, Host A (i.e., the iperf client) establishes a MPTCP connection with Host B (i.e., the iperf server). Afterwards, Host A sends traffic over the two subflows (SF1, and SF2) toward Host B, for a predetermined duration of 30 sec., trying to fill in the available link bandwidth; this type of traffic emulates bulk data transfer between two hosts. While this method would be considered equivalent to *timed transfers* (i.e., performance tests measuring the time needed to complete a 10MB tranfer between two end nodes), there have been some deviations reported in literature [4]. However, in this work we focus on *"length transfers"* [4], leaving "timed transfer" verification for future work.

### 3.3 Design of Experimentation Scenarios

The main idea behind the design of the experimentation scenarios has been the quest for a universally superior congestion control and packet scheduling algorithm combination, which could potentially be suitable for generic use, regardless of the underlying path conditions. And if the experimentation campaign would not come up with such a distinguishable combination, it would still be of interest to identify which packet scheduling and congestion control algorithm combination is the most appropriate under specific network conditions.

In an attempt to assess the various packet scheduling and congestion control algorithm combinations under various path heterogeneity conditions, the Mininet emulator has been used. Execution in Mininet is performed over the real protocol stack instead of a simulated protocol model. Furthermore, each test iteration has been executed using a static network configuration for the entire test duration (30 sec.), which in addition to the absence of background traffic, made our simple network setup ideal for assessing the combined use of algorithms, eliminating the effect of other factors on the performance outcome. This allowed us to repeat identical experiments and yet expect deterministic results, while providing us with the flexibility to modify the underlying path conditions and cover a wide range of path heterogeneity scenarios.

Our experimentation campaign included scenarios of variable path heterogeneity and severity in terms of capacity, round-trip delay, and packet loss rate, as well as all possible combinations of the packet scheduling and congestion control algorithms included in 1[11]and 2. This has been in accordance with the qualitative and quantitative factors that should be considered in such experiments, as mentioned in [72].

The entire set of experiments is depicted in Table 3.

In Table 3, SF1 and SF2 traffic characteristics are referred to the underlying Mininet link conditions (L1 and L2, respectively) over which the subflows have been established. Link L3 is statically configured to a $BW/RTT/PLR$ value of $2Gbps/0ms/0\%$ across the entire set of experiments. Based on the range of values described in Table 3 which determine the underlying path characteristics, we have initially divided our experimentation campaign in five

---

[11]ReMP TCP has been excluded from our experiments since, by duplicating packets, it clearly incurs sub-optimal performance in terms of goodput and overall flow completion time.

Table 3. Experimentation Scenarios

| Subflow | Capacity (Mbps) | Round-trip delay (ms) | Packet loss rate (%) |
|---------|-----------------|------------------------|----------------------|
| SF1 | [100] | [0, 5, 10] | [0, 1, 2] |
| SF2 | [100, 75, 50, 5] | [0, 2, 5, 20, 50] | [0, 0.05, 0.5, 2] |

major scenario families: *homogeneous*, *mild_heterogeneity*, *intense_heterogeneity*, *very_intense_heterogeneity*, and finally *mixed_heterogeneity*. Each scenario family, is composed of up to seven scenarios, where heterogeneity is assessed over individual traffic characteristics (i.e., only bandwidth-induced heterogeneity, while round-trip delay and packet loss maintain equal values among the two subflows), as well as over combinations of them. In all scenario families except the last one, one subflow (i.e., SF1) maintains static path configuration across all scenarios of the scenario family, while heterogeneity is induced by degrading path conditions of the second subflow (i.e., SF2 is inferior to SF1); in the scenarios belonging to the mixed_heterogeneity scenario family, both subflows experience mixed path characteristics, thus not permitting any a priori assumption as to which subflow is expected to perform better and how the applied algorithms would react (e.g., SF1 is established over a bottleneck link of $BW/RTT/PLR$ equal to $100Mbps/10ms/0\%$, while SF2 is over a $100Mbps/2ms/0.5\%$). These five scenario families produced twenty-nine (29) scenarios in total. Each scenario is executed against all combinations of five selected MPTCP packet schedulers (minRTT, BLEST, ECF, RR[12], LLHD) and seven congestion control algorithms (LIA, OLIA, BALIA, wVegas, Cubic, BBR, C-MPBBR); for each specific packet scheduling and congestion control algorithm combination, we execute the same scenario five times, to ensure that the outcome of each test is not affected by any transient issues. The total number of test iterations is given by the following formula:

$$\#\_tests = (\#\_scenarios) * (\#\_schedulers) * (\#\_CCAs) * (\#\_iterations\_per\_scenario) \tag{17}$$

which produces 5,075 test instances; as mentioned earlier, each test instance has a duration of 30 sec.

## 3.4 Performance Metrics

Inspired by the work performed in [13], we selected to provide a combined view of the per-subflow goodput, the per-subflow smoothed RTT (SRTT), the MPTCP out-of-order (OFO) queue size, as well as the per-subflow number of retransmissions for each packet scheduling algorithm.

*Goodput and Retransmissions.* Our performance evaluation results illustrate, among others, the achieved goodput which is essentially derived by subtracting the protocol overhead (i.e., packet headers) from throughput computation. While iperf tool provides a metric of the achieved goodput, this suffers two limitations: i) iperf tool provides the aggregate goodput achieved over both subflows of the MPTCP connection, while our intention was to depict how each packet scheduler performs, hence how traffic is actually distributed among the two subflows, and ii) iperf's goodput scale is calculated based on divisions by a factor of 1000. For these reasons, we have selected to analyze the per-subflow goodput, as well as to retrieve the per-subflow number of retransmissions by capturing and analyzing every single packet. The python-based pyshark tool [24] has been used for this purpose, allowing us to calculate goodput based on each individual packet's payload; our extracted per-subflow goodput scale has been calculated through divisions by a factor of 1024 (e.g., $1\,KiB = 1024\,bits$). For instance, iperf goodput expressed in Mbps using SI units (e.g., KB, MB)

---

[12]Round-robin packet scheduling algorithm maintained its default parameter values throughout the entire experimentation trial; that is, the number of consecutive segments sent on a subflow before switching to the other subflow is kept by default to one (num_segments = 1), and the scheduler tries to fill in the congestion window (cwnd_limited=true), which means that upon filling up the congestion window, becomes *ack-clocked* [6].

would be calculated by the formula: $GP(Mbps) = (Bytes/duration) * 8/(1000 * 1000)$, while our calculations based on binary units (e.g., KiB, MiB) are given by the formula: $GP(Mbps) = (Bytes/duration) * 8/(1024 * 1024)$, where $GP$ stands for *goodput*, and *duration* denotes the duration of each test instance (30 sec.). Thus, our goodput results might be even more conservative than those iperf produced. *Pyshark* allows the distinction between subflows via common tshark filters; this makes the calculation of per-subflow goodput and the respective number of retransmissions easy.

*SRTT and OFO queue.* The SRTT [81] as well as the OFO queue size values have been collected by tracing the respective linux kernel parameters. Noteworthy here is mentioning that for as long as the kernel tracing feature has been enabled (that is, for the 30 sec. duration of each test instance's execution), a huge amount of kernel logs containing metrics had been produced; thus, as part of a post-processing step, we randomly sampled the collected metrics in order to limit, for instance, the out-of-order queue value samples to two per millisecond, that is, 1 *sample* / 500*us*. Part of the collected metrics has also been the congestion window value of each subflow; however, we have selected to skip any congestion window plots and to focus on basic performance metrics.

Besides providing a direct performance comparison between the packet scheduling algorithms under a specific traffic scenario and under the presence of a single congestion control algorithm, the main focus of this work has been to assess the packet scheduling and congestion control algorithms across all traffic scenarios. For this reason, as it will be explained in detail within Section 4, a couple of additional assessment metrics have been used, such as the *normalized goodput* achieved by each packet scheduler across a scenario family; a *score* definition is also provided within Section 4, wherever deemed necessary. Finally, despite the fact that our traffic pattern resembles bulk transfer, we considered it useful to provide an overview of the per-packet delay which is an important metric in case of interactive applications. As such, we calculate the per-packet delay based on the following formulas:

    i. we first calculate the number of packets-per-second (pps):

$$pps = total\_Bytes/(MSS * test\_duration) \tag{18}$$

       where $MSS = 1514\,Bytes$ and test duration is $30\,sec$.

    ii. then, the per-packet-delay in milliseconds is given by:

$$per\_packet\_delay\,(ms) = (1/pps) * 1000 \tag{19}$$

Another performance metric to be also explained in detail within Section 4, is the *coefficient of variation (CV)* which is given by $\frac{\sigma}{\mu}$, and is used when calculating the congestion control algorithm scores as a measurement of the extent of variability in relation to the mean.

## 4  Evaluation

This section is dedicated to the overall evaluation of the packet scheduling and congestion control algorithms. We first present a subset of figures, illustrating how the MPTCP packet scheduling algorithms performed across different congestion control algorithm (CCA) combinations, under some representative scenarios which resemble real-world network conditions. Then, we define two new metrics to assess the performance of the packet scheduling and the congestion control algorithms, that is, the *PS_score* and the *CCA_score*, respectively. The *PS_score* metric is based on the normalized goodput in order to assess the performance level of each packet scheduling algorithm within each scenario family. The *CCA_score* metric, which is based on the normalized *PS_score*, is used to assess the performance of the congestion control algorithms across all sccenario families. This metric provides the means to quantify the

performance of each CCA through the normalized score that the underlying packet scheduling algorithms achieved using the respective CCA. The *CCA_score* indicates, in essence, how much "space" each CCA provides to the underlying packet scheduling algorithms to unveil their sophisticated scheduling logic. If a congestion control algorithm suppresses the congestion window of a subflow, then the packet scheduling algorithm cannot perform any better than applying its logic on the resources that the CCA makes available to the scheduler; this means that a bad congestion control algorithmic logic is anticipated to incur a poor performance outcome, regardless of the packet scheduling algorithm's intrinsic properties.

Once we identify the "best" three congestion control algorithms, we proceed with a fine-grained assessment of the various packet scheduling algorithms when combined with each of the highest-score CCAs. The *normalized goodput* and the *per-packet delay* are the main metrics against which all packet scheduling algorithms are evaluated. The `Matplotlib` python library has been used for the visualization of the results [45].

Noteworthy here is that we have deliberately decided to maintain different axes' scale when plotting the values of some parameters across the various CCAs, such as the SRTT, the MPTCP OFO queue size, and the number of retransmissions. This allows to increase the granularity and, consequently, the visibility of the results, preventing subtle variations from being shadowed by any neighboring CCA's more coarse results within the same figure.

### 4.1 Cross-CCA packet scheduling overview

From the 1,015 unique tests executed (5,075 tests / 5 iterations per test), we have selected to present in this work a subset of the most representative scenarios. To this end, we hereby present:

    i. a scenario involving homogeneous paths,
    ii. a scenario involving paths experiencing intense heterogeneity, and
    iii. a scenario involving path conditions of mixed heterogeneity.

*Homogeneous scenario.* Figure 5 presents the results of a scenario derived from the *homogeneous* scenario family, where both subflows traverse paths which experience common networking characteristics. When an equal delay of 5ms is present on both paths, minRTT, ECF, and RR outperform LLHD and BLEST under LIA CC, and its variations OLIA and BALIA. All schedulers follow the same performance pattern across these three CCs. It is evident that even when the CCA (i.e., LIA/OLIA/BALIA) leaves the congestion window of both subflows open, BLEST, which bases its decisions on send window blocking estimations, selects to equally limit the use of both subflows, underutilizing the available bandwidth of both paths, and eventually underperforming. wVegas, which is a delay-based CC, is highly sensitive to delay, interpreting it to congestion and capping aggressively the congestion window of both subflows; the aggregate goodput achieved by all schedulers is limited compared to the one achieved under the use of other CCAs. On the contrary, the use of Cubic, which is a loss-based uncoupled CCA, provides a clear benefit to all schedulers. BLEST's performance improvements can be credited to the uncoupled nature of Cubic, that is, there is no upper bound for the aggregate congestion window of the two subflows, as is the case for coupled CCAs. Thus, each of the subflows maintains its individual congestion window which is treated by the employed CCA as if it were for a singe path flow belonging to a normal TCP connection. This makes MPTCP unfair to single path TCP flows, since each MPTCP subflow can allocate a portion of BW equal to that of a single path TCP flow, ultimately leading single path TCP flows to starvation when single path TCP and MPTCP subflows coexist in the same bottleneck link. In this sub-scenario, BBR allows all schedulers to reach high and stable performance; it is also noticeable that BBR's accurate estimations of the bottleneck bandwidth and RTT lead to scheduling decisions which incur the lowest SRTT and very low OFO queue occupancy across all the

employed CCAs. C-MPBBR, which is a coupled variant of BBR for MPTCP (ensuring that MPTCP has a clear benefit over and is fair to SP-TCP), incurs a lower goodput aggregation to all schedulers' performance, when compared to uncoupled BBR. The fairness C-MPBBR achieves across the two MPTCP subflows traversing the same bottleneck link, can be clearly seen in the measured SRTT of the subflows for each scheduler; however, its coupled nature and fairness benefits come at the cost of a more restricted aggregate congestion window which limits the performance of all schedulers over the two subflows. The overall $e1\_hom\_bw\_delay$ scenario's verdict can be summarized as follows:

- minRTT, ECF, and RR exhibit high performance under LIA/OLIA/BALIA CCAs
- BLEST is sensitive to delay and underperforms the rest when any one of the LIA/OLIA/BALIA is the employed CCA
- wVegas is a delay-based CCA and reacts more aggressively to delay, restricting the available congestion window, and leading all schedulers to low goodput aggregation over the two paths
- all schedulers perform comparably well under Cubic
- all schedulers perform comparably well and stably under BBR
- C-MPBBR bounds the coupled congestion window in favor of fairness and MPTCP's benefit over SPTCP; this has a considerable impact on the performance of all schedulers when compared to that achieved with C-MPBBR's rivals (i.e., uncoupled BBR and Cubic)
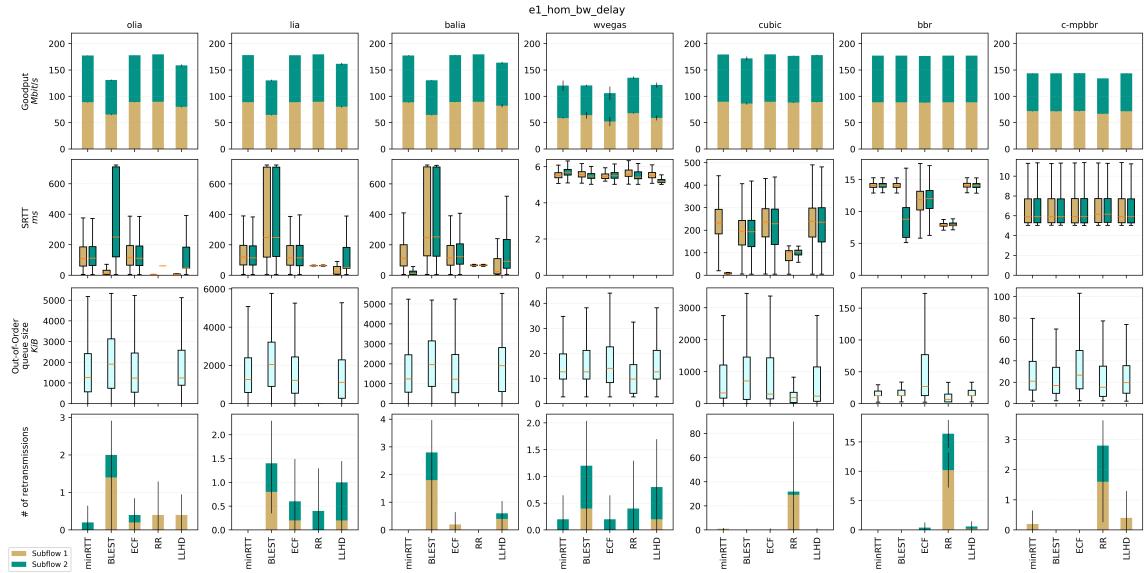


Fig. 5. Homogeneous scenario - BW/RTT/PLR - $SF_1$ : $100\,Mbps/5\,ms/0\%$, $SF_2$ : $100\,Mbps/5\,ms/0\%$

*Intense_heterogeneity scenario.* Figure 6 presents the results of a scenario derived from the *intense_heterogeneity* scenario family, where the two subflows traverse paths with intense-heterogeneity networking characteristics. The additional presence of delay on a path (i.e., SF2) on top of intense bandwidth restrictions and packet loss, has a detrimental effect on the goodput of all schedulers, when combined with any one of the traditional loss- and delay-based CCAs. All schedulers fail to efficiently utilize the available bandwidth of the affected path (SF2) when combined with

any one of the LIA/OLIA/BALIA/wVegas/Cubic; this result is mainly credited to the effects of the congestion control
approaches that the aforementioned CCAs follow; that is, increasing the congestion window during the CA (congestion
avoidance) phase, which in turn allows intermediate buffers to be filled up, and then trying to tackle the congestion
they created by drastically decreasing the congestion window, which deteriorates the already adverse conditions even
further. Under such circumstances, schedulers cannot perform any better than trying to efficiently utilize the resources
that the CCA makes available; CCA has a great impact (either positive or negative) on the final performance. The
craftsmanship in the design of BBR and C-MPBBR, lies in the way they estimate the bottleneck bandwidth and the
round-trip propagation time, thus pacing traffic (via the cwnd) at the highest possible rate which, however, does not lead
to buffer formation; when buffers start being occupied, BBR can detect the RTT increase and adjust the rate efficiently
until the queue is drained and a lower RTT is measured. This different logic allows BBR to sustain traffic rate close
to Kleinrock's optimal operating point, avoiding the post-congestion control operation of traditional CCAs (AIMD).
BBR and C-MPBBR are rather adjusting traffic smoothly by avoiding high queue occupancies and the subsequently
increased congestion; this allows all schedulers to unleash their full potential and apply their scheduling logic towards
high goodput rates. In this sub-scenario, minRTT, ECF, RR, and LLHD perform comparably well, achieving close to
optimal goodput; BLEST is not limited by the CCA decisions, but rather applies its algorithmic logic, which is based on
limiting or eliminating the use of a slow path if this is anticipated to cause HoL and, consequently, MPTCP send window
blocking (due to inflight data carried over the slower path). According to this logic, BLEST selects not to use the affected
path (SF2), underutilizing its available bandwidth and thus exhibiting the lowest performance among schedulers. The
overall *e7_int_het_loss_bw_delay* scenario's verdict can be summarized as follows:

- the decisions of traditional delay- and loss-based CCAs (i.e., LIA/OLIA/BALIA/wVegas/Cubic) have a severe
  impact on any scheduler's performance; when these CCAs have to deal with intense delay in conjunction with
  capped bandwidth and packet loss, their decisions have a disastrous effect on the final goodput, irrespective of
  the employed scheduler
- BBR (and its MPTCP coupled variant C-MPBBR) are designed in a way that does not aggressively suppress the
  congestion window of an already adversely impacted flow, thus allowing schedulers to achieve and sustain
  high goodput
- BLEST avoids the use of the slow path, either because the congestion window of the respective subflow is
  suppressed by a traditional CCA (i.e., LIA/OLIA/BALIA/wVegas/Cubic), or because of its own algorithmic logic
  to avoid slow paths in case this is expected to cause send window blocking
- RR experiences the highest number of retransmissions among schedulers

*Mixed_heterogeneity scenario.* Figure 7 presents the results of a scenario derived from the *mixed_heterogeneity*
scenario family, where the two subflows traverse paths experiencing mixed-heterogeneity networking conditions[13]. In
this particular mixed delay- and loss-induced heterogeneity sub-scenario, there is a clear distinction between schedulers'
achieved goodput across CCAs, illustrating the effect that the various CCAs may have in scheduling. The highest
scheduling performance is achieved when schedulers are combined with BBR; in that case, all schedulers (except for
RR) achieve comparably high goodput over both subflows. Equally high goodput rates (to those observed with BBR) are
reached by schedulers in case of uncoupled loss-based CCA Cubic; RR however is severely degraded. The next higher

---

[13]Mixed heterogeneity scenarios emulate complex networking conditions where each of the paths is superior to the other when considering only a single
factor; be it the allocated bandwidth, or the experienced RTT, or the packet loss rate. In such cases, however, it is sometimes not evident beforehand
which path is eventually going to render an overall better performance.
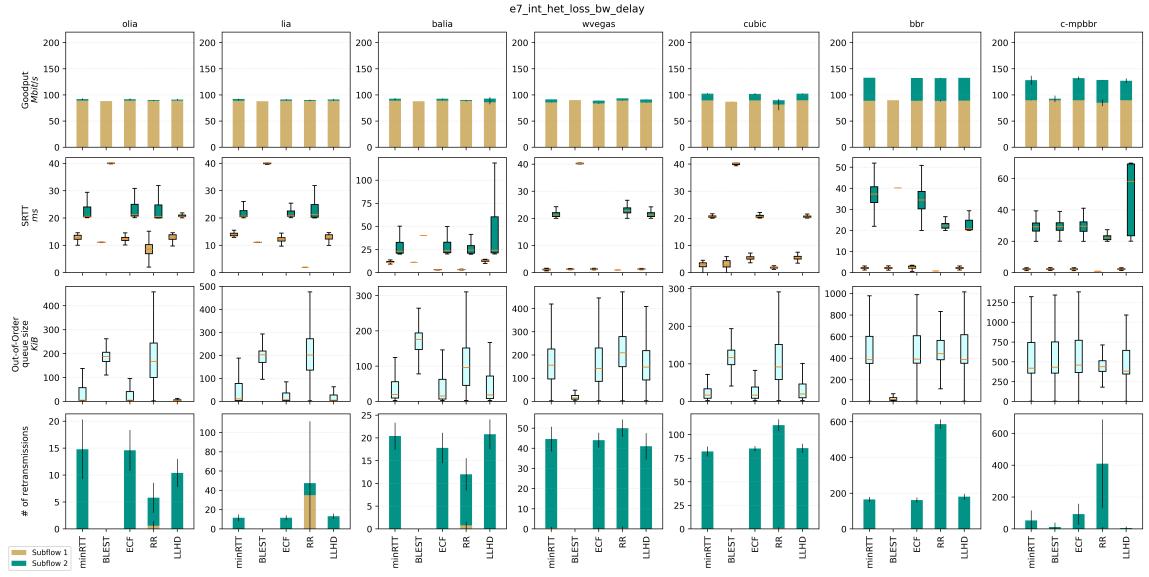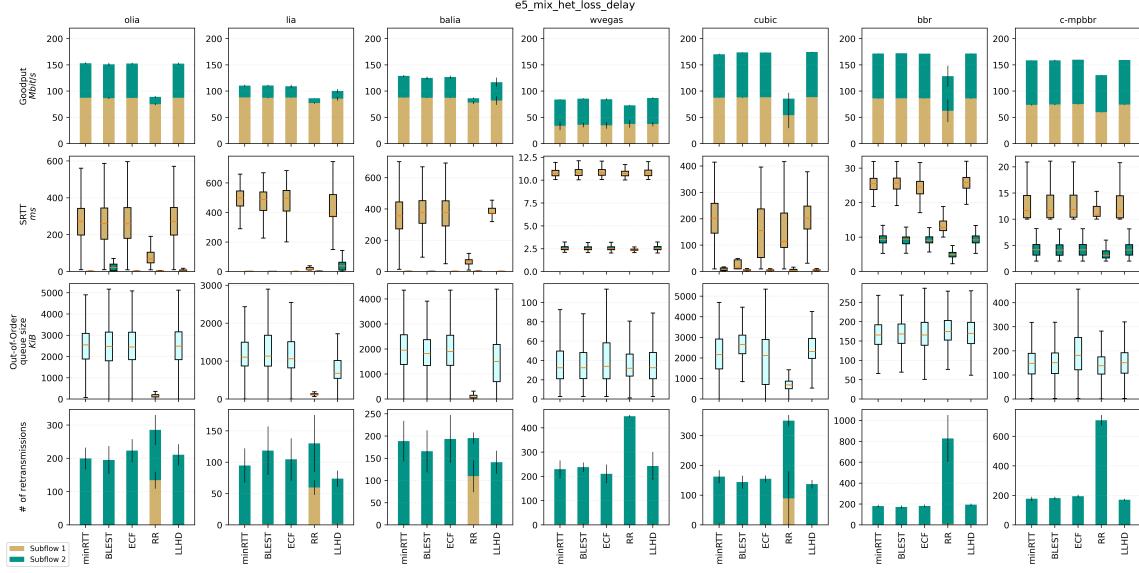
Fig. 6. Intense heterogeneity scenario - BW/RTT/PLR - $SF_1 : 100\,Mbps/0\,ms/0\%$, $SF_2 : 50\,Mbps/20\,ms/0.5\%$

goodput rates are achieved when schedulers are combined with C-MPBBR; all (except RR) are performing comparably high, but they experience some goodput restriction over the higher-delay path (SF1), as compared to that achieved with BBR. This can be credited to the restrictions that C-MPBBR's coupled nature imposes on the subflows that are considered to be traversing the same bottleneck path. With the couped loss-based OLIA, the goodput achieved over SF2 is slightly restricted (compared to the one achieved with BBR / C-MPBBR / Cubic), since the coupled nature of OLIA in conjunction with its more aggressive reaction to the packet loss detected over SF2, lead to further cwnd restrictions. The same goodput pattern among schedulers, with a higher goodput restriction over SF2, is also visible with BALIA, while with LIA, SF2's goodput is limited even further. Finally, wVegas has a detrimental effect on schedulers' performance; all schedulers underperform over both subflows, being able to eventually reach an aggregate goodput which is less than a half of the one achieved with BBR. The overall *e5_mix_het_loss_delay* scenario's verdict can be summarized as follows:

- minRTT, BLEST, ECF, and LLHD perform comparably to one another, within each CCA combination
- BBR, C-MPBBR, and Cubic provide ideal conditions to the underlying schedulers, allowing them to achieve very high goodput rates
- OLIA is the third best (in terms of schedulers' performance) CCA, allowing all schedulers to reach high goodput rates; goodput over the lossy path (SF2) is restricted more due to the loss-based nature of OLIA which caps the cwnd to tackle congestion, not being able to distinguish random loss from actual link congestion
- BALIA and LIA limit the SF2 cwnd further, leading all schedulers to moderate performance
- wVegas has a detrimental effect upon all schedulers' performance, due to its delay-based congestion inference
- RR experiences the highest number of retransmissions and underperforms, regardless of CCA selection

Fig. 7. Mixed heterogeneity scenario - BW/RTT/PLR - $SF_1 : 100\,Mbps/10\,ms/0\%$, $SF_2 : 100\,Mbps/2\,ms/0.5\%$

## 4.2 PS and CCA score metrics

As mentioned earlier, to evaluate the performance of the various packet scheduling and congestion control algorithms, we have introduced two new metrics; the *PS_score* and the *CCA_score*, respectively.

*PS_score.* For each scheduler, we first calculate the *average goodput* it achieved in each sub-scenario of each broader scenario family, over the five identical execution runs. Then, we calculate a *PS_score* for each packet scheduler per scenario family; the *PS_score* is actually the *normalized aggregate goodput* for each scheduler in each scenario family, and is given by the formula:

$$PS\_score = \sum_{\#subscen} \left( \frac{subscen\_avg\_GP}{\#subscen\_in\_ScenFam * theoretical\_max\_gp} \right) \tag{20}$$

where *#subscen* denotes the number of sub-scenarios included in each respective scenario family, *subscen_avg_GP* is the average goodput achieved over the five identical iterations of each respective sub-scenario, *#subscen_in_ScenFam* denotes the number of sub-scenarios included within a broader scenario family, and the *theoretical_max_gp* denotes the theoretical maximum aggregate goodput over the two subflows that is, $200\,Mbps$. This way, $PS\_score \in [0, 1)$.

*CCA_score.* A similar approach has been followed in order to define a metric which would allow us to assess the performance of the congestion control algorithms. We first present how each congestion control algorithm performed within each broader scenario family. This kind of assessment is essentially showcasing how the packet scheduling algorithms performed under each respective congestion control algorithm within each broader scenario family. Thus, each CCA's *CCA_score_per_ScenFam* is computed based on the average *PS_score* of the packet scheduling algorithms

when combined with the respective CCA. The *CCA_score_per_ScenFam* is provided by the following formula:

$$CCA\_score\_per\_ScenFam = \sum_{\#schedulers} \left( \frac{PS\_score}{\#schedulers} \right) \tag{21}$$

where *#schedulers* denotes the number of all packet scheduling algorithms, and *PS_score* is the respective score each packet scheduling algorithm received for this scenario family. For this metric, it also applies that *CCA_score_per_ScenFam* ∈ [0, 1).

### 4.3   Congestion control assessment

The *CCA_score_per_ScenFam* reached by each congestion control algorithm in each scenario family, is depicted in Figure 8.
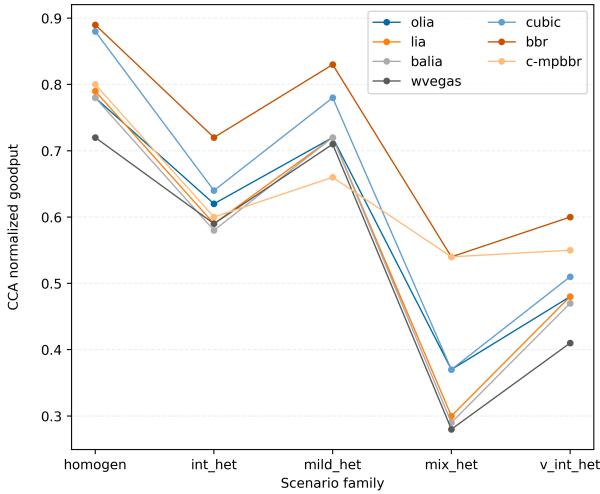


Fig. 8.  CCA score per scenario family

Once an overview had been obtained on how each congestion control algorithm performed within each scenario family, it was interesting to identify the three best-performing CCAs across all scenario families. To this end, we defined the *CCA_score* which is essentially an average over the scores each congestion control algorithm has been assigned for its performance within each scenario family. *CCA_score* is consequently a *cross-scenario-family* score indication, given by the formula:

$$CCA\_score = \sum_{\#ScenFam} \left( \frac{CCA\_score\_per\_ScenFam}{\#ScenFam} \right) \tag{22}$$

where *#ScenFam* is the number of the different scenario families.

The drawback, however, of the above-mentioned metric is the fact that it does not consider how poorly each congestion control algorithm performed in its worst-case scenario family. Thus, to allow also for performance stability across the various scenario families, we combined the *CCA_score* with its respective *coefficient of variation (CV)*, which provides insight on the extent of variability in relation to the mean. Figure 9 illustrates the coefficient of variation for each individual *CCA_score*.

The lower the coefficient of variation, the higher the performance stability of each CCA across the various sub-scenarios. Dividing each CCA's *CCA_score* by its respective coefficient-of-variation value, we derive an *overall score* for each congestion control algorithm, which allowed us to identify the best three in terms of goodput and performance stability across scenario families. The *CCA_overall_score* is then given by the below formula:

$$CCA\_overall\_score = \frac{CCA\_score}{CCA\_CV} \tag{23}$$

Figure 10 illustrates each congestion control algorithm's *CCA_overall_score*. This way, we have been then able to identify the three best-performing congestion control algorithms: **C-MPBBR**, **BBR**, and **OLIA**.
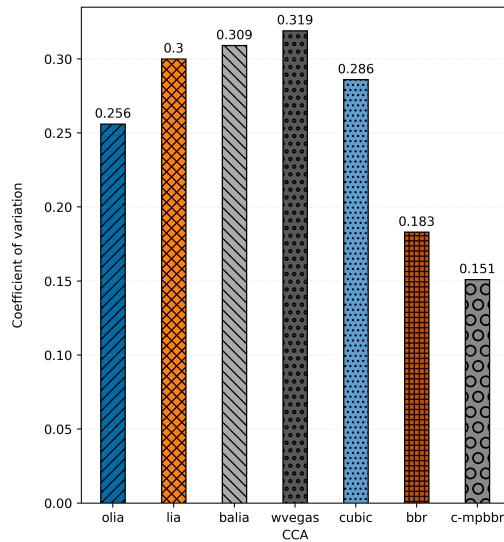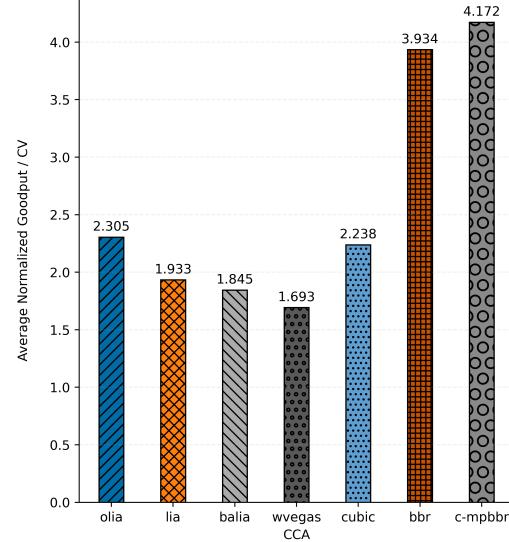


Fig. 9.  CCA_score CV



Fig. 10.  CCA_overall_score

## 4.4 Packet scheduling assessment - Goodput

Having identified the three best-performing CCAs, we have then been able to delve into the packet scheduling algorithms' performance evaluation. The *PS_score* metric is indicative on how each scheduler performed within each scenario family, under each of the best three congestion control algorithms. Figure 11 illustrates each packet scheduling algorithm's performance score per CCA, and per scenario family.

A more granular *PS_score* analysis per scenario family is provided for each one of the BBR, C-MPBBR, and OLIA CCAs, within Figures 12, 13, and 14, respectively.

An alternative view of each packet scheduling algorithm's score per scenario family for each of the top three CCAs, is provided via the *PS_score* heatmaps showcased within Figures 15, 16, and 17, respectively.

One of the main purposes of this work, has been to provide a clear view on how each packet scheduler performs across all scenario families for each of the best-performing CCA. To this end, Figure 18 illustrates the *Empirical Complementary Cumulative Distribution Function (ECCDF),* which depicts each scheduler's probability to perform above a specific level
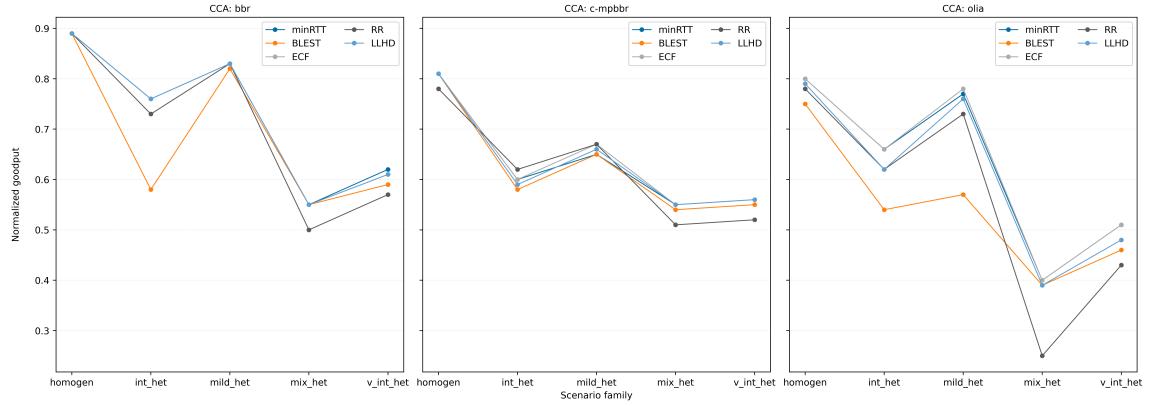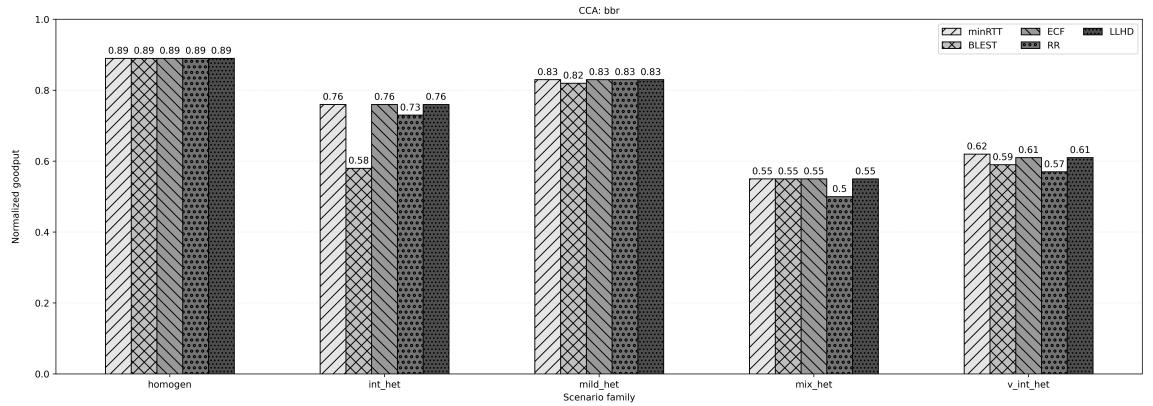
Fig. 11. PS_score



Fig. 12. PS_score per scenario family for BBR

of goodput, for each of the top-three CCAs (i.e., BBR, C-MPBBR, and OLIA); ECCDF considers goodput probability across all scenario families and sub-scenarios.

## 4.5 Packet scheduling assessment - Per-Packet Delay

Besides the evaluation in terms of goodput, we found it also useful to provide insight on how each packet scheduling algorithm would perform in case of interactive applications. A key metric for this type of applications is the experienced *per-packet delay (PPD)*. For this purpose, we provide a high-level view of each packet scheduling algorithm's[14] performance in terms of per-packet delay in Figure 19; a more fine-grained numerical evaluation of each packet scheduling algorithms per-packet delay for each of the best-three congestion control algorithms (i.e., BBR, C-MPBBR, and OLIA) can be found in Figures 20, 21, and 22, respectively.

---

[14]"Packet Scheduling Algorithm" is abbreviated as *PSA* within figure captions.

Fig. 13. PS_score per scenario family for C-MPBBR



Fig. 14. PS_score per scenario family for OLIA



Fig. 15. PS_score heatmap per scenario family for BBR

Fig. 16. PS_score heatmap per scenario family for C-MPBBR

Fig. 17. PS_score heatmap per scenario family for OLIA

Finally, Figure 23 depicts the *Empirical Cumulative Distribution function (ECDF)*, which provides a metric of scheduler's probability to achieve below a specific level of average per-packet delay, for each one of the top-three CCAs (i.e., BBR, C-MPBBR, and OLIA); ECDF considers the *average per-packet delay* across all scenario families and sub-scenarios.

Fig. 18.  ECCDF depicting packet scheduler goodput probability across scenario families, for BBR, C-MPBBR, and OLIA



Fig. 19.  Qualitative view of each PSA's PPD per scenario family for BBR, C-MPBBR, and OLIA



Fig. 20.  PPD of each PSA per scenario family for BBR

Fig. 21.  PPD of each PSA per scenario family for C-MPBBR
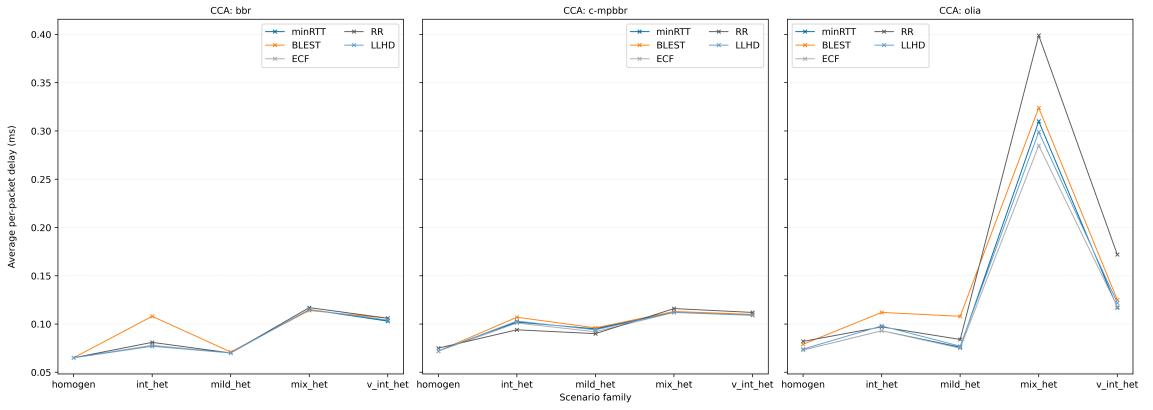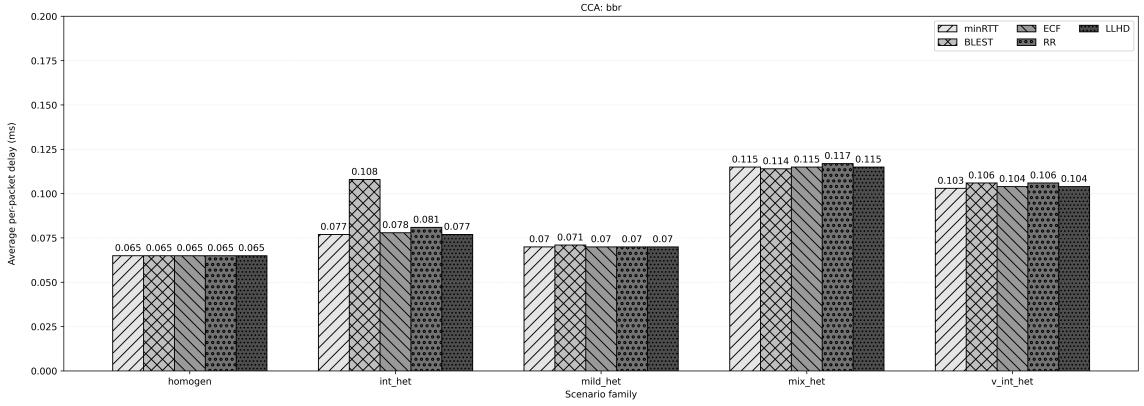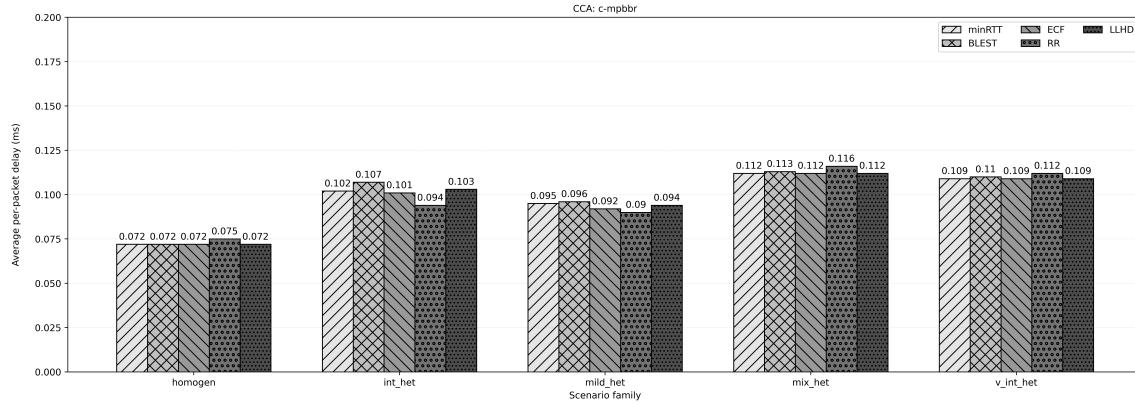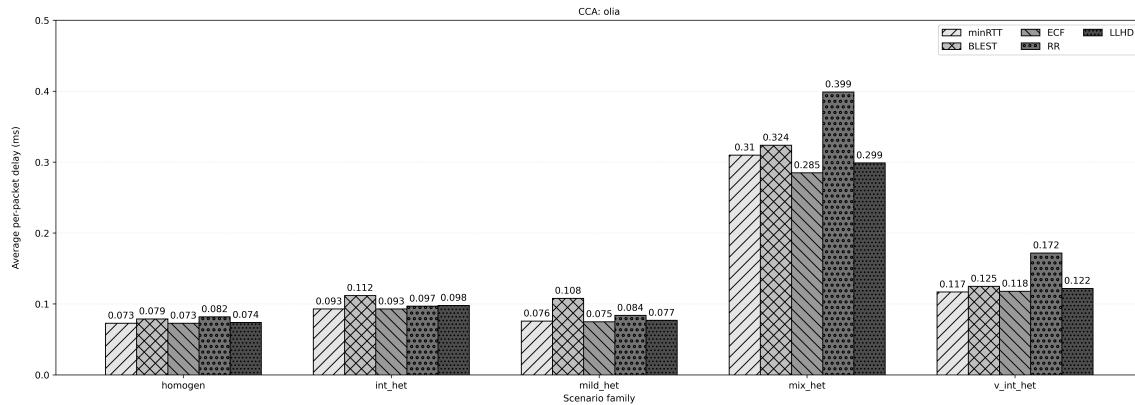


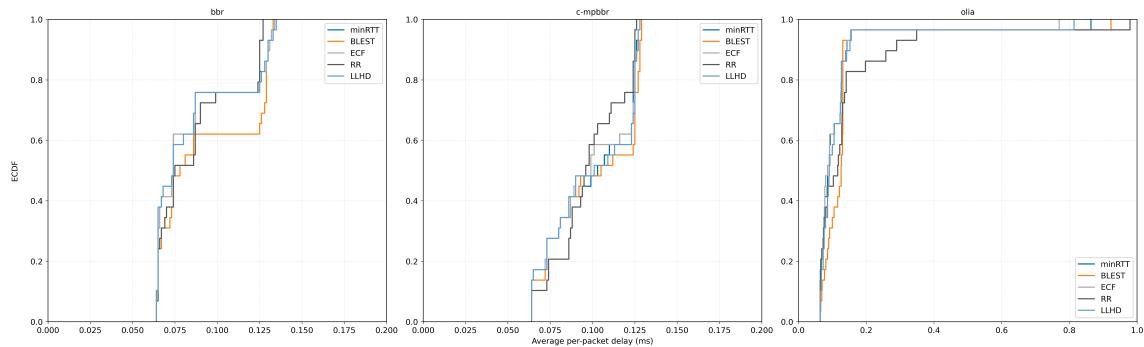Fig. 22.  PPD of each PSA per scenario family for OLIA



Fig. 23.  ECDF depicting packet schedulers' PPD probability across scenario families, for BBR, C-MPBBR, and OLIA.

## 5   Discussion

The performance results derived by this extended experimentation campaign indicate that the there is no single congestion control and packet scheduling algorithm combination that fits well all underlying network conditions. As a result, it would be recommended to select an algorithmic scheme combination according to the level of heterogeneity experienced in the network, as illustrated through the aforementioned results. On the other hand, adjusting the protocol configuration for fine-tuning can be a cumbersome procedure since network changes can be of dynamic, and in most cases stochastic, nature. It is then reasonable to seek for a stable configuration, able to cope with all potential network characteristics. Toward this objective, the initial factor which determines the overall performance outcome, is the appropriate selection of the congestion control algorithm. In this direction, our analysis indicated that *BBR*, *C-MPBBR*, and *OLIA* could be potentially applied as generic-use congestion control schemes. With any one of these three CCA's, our performance results indicate that *ECF* and the *MPTCP default (minRTT)* packet scheduling algorithms exhibit the highest goodput and lowest per-packet delay performance, across almost all scenarios. *RR* turns out to be useful only for experimentation purposes, and should rather be avoided in production environments, since it experiences the worst results among the examined schedulers; RR's performance outcome across the diverse scenarios confirm the recommendation included within [6]. *LLHD* is also considered a stable, in terms of goodput and per-packet delay, packet scheduling algorithm whose results are very close to these of the first two. *BLEST* is also a highly-sophisticated algorithm which performs quite well in case of homogeneous and mild heterogeneity; however, when the underlying network's path heterogeneity becomes severe, BLEST decides to skip the non-beneficial paths, underutilizing the available network resources. Doing so, diminishes its performance results compared to the other best-performing schedulers (except RR which performs the worst). Nevertheless, even with this kind of resource underuse, BLEST does not experience any significant performance inefficiency, and can thus be included within the group of top-four. An important observation here is the fact that if the optimization goal is application-level smoothness instead of raw throughput maximization, BLEST could be considered superior since it actively avoids the impaired path to prevent HoL blocking. This attribute also stresses the fact that the "best" scheduler is ultimately application-dependent.

One limitation of the current experimentation test-set is the fact that our traffic resembled bulk-traffic conditions; thus, schedulers that might lag behind in this specific use case, might reach better results under other types of traffic, such as short web-page loads, or traffic generated by video streaming applications.

## 6   Related Work

There are different directions in literature around the performance assessment of MPTCP and its underlying packet scheduling and congestion control algorithms. Some works maintain the default MPTCP algorithm configuration (e.g., minRTT packet scheduler and (O)LIA congestion control) and examine the impact of the various (MP)TCP protocol parameters (e.g., send/receive buffer size, SACK, TCP auto-tuning, etc.) on the overall MPTCP performance. Another direction is focusing on a specific research area (e.g., on packet scheduling) and assessing how the various algorithms of this particular domain perform, while keeping the default configuration in other protocol areas (e.g., default congestion control and/or default (MP)TCP parameters). Finally, there is a limited number of works assessing how the various algorithm combinations affect MPTCP protocol operation. Our work is positioned to this latter case, coupling the theoretical aspects with an extensive performance assessment of the state-of-art packet scheduling and congestion control algorithms. Our manuscript's motivation and distinctive aspects are described hereinafter.

The performance evaluation conducted within [1] has been one of the first attempts to quantify the impact that various factors may impose on MPTCP performance (i.e., the buffer size, packet loss, the MSS-size). Besides evaluating the impact of a single parameter (i.e., BW, RTT delay, or packet loss) on MPTCP performance, our work considers also a mixture of these parameters and examines how these affect factors other than merely goodput, such as the SRTT, the OFO-queue size, and the number of retransmissions. Moreover, our experiments go beyond the default MPTCP packet scheduling (minRTT), and Reno / Coupled (LIA) congestion control schemes, assessing a variety of schedulers and congestion control algorithm combinations. [71] implements the first framework which enables loading packet scheduling and congestion control algorithms as modules, allowing algorithm alternations at runtime. Experimental evaluation there, is conducted using minRTT scheduler, both for bulk transfer and application-limited traffic. On the contrary, our work examines a bunch of modern packet scheduling and congestion control schemes. [43] is a comprehensive work trying to measure the extent to which MPTCP is beneficial compared to single-path TCP under different scenarios. The authors consider minRTT scheduler combined with Coupled (LIA) CCA for MPTCP, and Cubic CCA for SP-TCP, while our work examines multiple packet scheduling and congestion control algorithm combinations. Another difference is the fact that the experiments conducted throughout [43] consider variable traffic size and higher 5G bandwidths close to 1Gbps; our work seeks to provide foundational details behind the operation of schedulers and CCA combinations, thus selected to keep a static length-based traffic approach (iperf traffic tries to fill in resources for 30 sec. test duration) and cap bandwidth to 100Mbps per subflow in order to be able to identify basic discrepancies and inefficiencies within the algorithms, which might be invisible in high data rates. [46] focuses on the development of a new scheduler (STTF) while assessing it against other packet scheduling algorithms (i.e., minRTT, BLEST, ECF, DAPS [61], OTIAS [91]) under various traffic conditions. However, packet schedulers are compared to one another using only the default MPTCP congestion control algorithm, which was Coupled (LIA). [76] provides a high level overview on the advancements of transport layer protocols, a classification of the various congestion control schemes and a short description of them, as well as details around the various multipath protocols. In a different direction, our work is focused on a limited number of packet scheduling and congestion control schemes, providing an in-depth analysis of their internal operation as well as the detailed algorithmic steps. Furthermore [76] is a theoretical survey, while our work is an experimental manuscript, combining theory with experimental evaluation of the algorithms under review. As mentioned earlier, part of our work has been inspired by the work conducted in [13], especially regarding the visibility of the results in a combined figure including goodput, SRTT, OFO queue size, and the number of retransmissions. [13] also provides a short description as well as the algorithmic steps of some packet scheduling algorithms (i.e., minRTT, RR, MuSher [80], BLEST, ECF, STTF [46]), and examines their performance in conjunction with some congestion control schemes (i.e., LIA, OLIA, BALIA, wVegas). However, it modifies the send and receive buffers to identify their impact on performance, while we selected to maintain the default values in order to focus mainly on the underlying algorithms' performance, and to avoid blending many factors which could further affect performance results. In addition, our work includes the latest BBR and C-MPBBR congestion control schemes and tries to further analyze the results by introducing novel metrics, in an attempt to directly assess all possible combinations and extract the best-performing ones. Another extensive experimental survey has been provided through the work conducted in [60]. The authors of this work provide a comprehensive analysis of packet scheduling schemes along with a classification based on their use cases and the number of scheduling criteria considered by each algorithm. Their work has also influenced the design of our own experimentation strategy, and considers a nearly-exhaustive list of packet schedulers and CCAs. However, the authors of that work examine different packet scheduling and a subset of the congestion control algorithms we use in our own experiments (i.e., $PSAs = \{minRTT, RR, LWS, LTS, HSR\}$ and $CCAs = \{LIA, BALIA, OLIA, wVegas\}$, respectively).

## 7  Conclusion

This work has been conceived during our own research on multipath protocols and more specifically on MPTCP. While anyone interested to conduct research in this field may need to skim through multiple sources, we thought it would be useful for the reader to have a single-point of reference where all the basic information around MPTCP protocol is collected. To this end, Section 2 has been dedicated to providing the necessary details for a preliminary introduction to the MPTCP protocol, its structure and architectural overview, as well as the main modules determining its operation. This section provides also detailed information on the various packet scheduling and congestion control schemes, including both their theoretical and algorithmic aspects. To this end, five packet scheduling (i.e., minRTT, BLEST, ECF, RR, and LLHD) and seven congestion control algorithms (i.e., Cubic, Coupled/LIA, OLIA, BALIA, wVegas, BBR, and C-MPBBR) have been examined, in an attempt to cover the theoretical aspects before moving on with their performance evaluation. Section 3 provided the overall experimentation methodology, our Mininet environment topology, as well as the underlying protocol and kernel configuration details, which we considered when preparing our experimentation testbed. The same section included also details of the traffic methodology, the design of the experimentation scenarios, and the performance metrics used to assess the various algorithm combinations. Section 4 presented the experimentation results of the various packet scheduling and congestion control algorithm combinations. The score metrics conceived to assist our evaluation as well as the main evaluation criteria (e.g., goodput and per-packet delay), have been also presented within this section. Section 5 concluded the core part of this work by providing a short summary on the best-performing packet scheduling and congestion control schemes, while Section 6 provided details on the related work present in literature.

While this work serves as a starting point for anyone interested in studying the preliminaries of MPTCP, such as the basic details around the protocol operation, as well as the available packet scheduling and congestion control schemes, it is not and cannot be a single point of study; however, appropriate references are provided to the original manuscripts and IETF RFCs, which provide the entire picture.

Finally, the list of packet scheduling and congestion control algorithms that have been selected for our theoretical analysis and the experiments, contains the ones considered as foundational and best-performing for MPTCP, laying essentially the groundwork for the development of many other algorithms, either unique or variations of existing ones.

## Acknowledgments

## References

[1]  Sébastien Barré et al. 2011. *Implementation and assessment of modern host-based multipath solutions*. Ph. D. Dissertation. Catholic University of Louvain, Louvain-la-Neuve, Belgium.

[2]  Ethan Blanton, Dr. Vern Paxson, and Mark Allman. 2009. TCP Congestion Control. RFC 5681. doi:10.17487/RFC5681

[3]  Olivier Bonaventure, Christoph Paasch, and Gregory Detal. 2017. Use Cases and Operational Experience with Multipath TCP. RFC 8041. doi:10.17487/RFC8041

[4]  Lawrence Sivert Brakmo. 1996. *End-to-end congestion detection and avoidance in wide area networks*. The University of Arizona.

[5]  Lawrence S. Brakmo, Sean W. O'Malley, and Larry L. Peterson. 1994. TCP Vegas: new techniques for congestion detection and avoidance. *SIGCOMM Comput. Commun. Rev.* 24, 4 (Oct. 1994), 24–35. doi:10.1145/190809.190317

[6]  C. Paasch, S. Barre, et al. 2023. *Multipath TCP in the Linux Kernel*. Retrieved August 28, 2025 from https://www.multipath-tcp.org

[7]  Canonical Ubuntu. 2024. *Download Ubuntu Desktop*. Canonical Ltd., London, UK. https://ubuntu.com/download/desktop

[8]  Y. Cao and E. Dong. 2021. wVegas CCA - Multipath TCP Linux Kernel v0.96. https://github.com/multipath-tcp/mptcp/blob/mptcp_v0.96/net/mptcp/mptcp_wvegas.c

[9] Yu Cao, Mingwei Xu, and Xiaoming Fu. 2012. Delay-based congestion control for multipath TCP. In *2012 20th IEEE international conference on network protocols (ICNP)*. IEEE, 1–10.

[10] Neal Cardwell, Yuchung Cheng, C Stephen Gunn, Soheil Hassas Yeganeh, and Van Jacobson. 2016. Bbr: Congestion-based congestion control: Measuring bottleneck bandwidth and round-trip propagation time. *Queue* 14, 5 (2016), 20–53.

[11] Neal Cardwell, Yuchung Cheng, Soheil Hassas Yeganeh, and Van Jacobson. 2017. *BBR Congestion Control*. Internet-Draft draft-cardwell-iccrg-bbr-congestion-control-00. Internet Engineering Task Force. https://datatracker.ietf.org/doc/draft-cardwell-iccrg-bbr-congestion-control/00/ Work in Progress.

[12] Neal Cardwell, Yuchung Cheng, Soheil Hassas Yeganeh, Ian Swett, and Van Jacobson. 2022. *BBR Congestion Control*. Internet-Draft draft-cardwell-iccrg-bbr-congestion-control-02. Internet Engineering Task Force. https://datatracker.ietf.org/doc/draft-cardwell-iccrg-bbr-congestion-control/02/ Work in Progress.

[13] Hendrik Cech. 2020. Analyzing and realizing multipath TCP schedulers in linux. *Technical University of Munich. PDF-dokumentti. Saatavissa: https://www. nitindermohan. com/documents/student-thesis/HendrikCechGR. pdf [viitattu 11.3. 2025]* (2020).

[14] Yuchung Cheng, Neal Cardwell, Soheil Hassas Yeganeh, and Van Jacobson. 2022. *Delivery Rate Estimation*. Internet-Draft draft-cheng-iccrg-delivery-rate-estimation-02. Internet Engineering Task Force. https://datatracker.ietf.org/doc/draft-cheng-iccrg-delivery-rate-estimation/02/ Work in Progress.

[15] Ask Ubuntu community. 2019. Enable TCP probe. https://askubuntu.com/questions/1164562/tcp-probe-could-not-insert-tcp-probe

[16] GitHub BBR community. 2024. Google BBR. https://github.com/google/bbr

[17] MPTCP community. 2021. minRTT - Multipath TCP Linux Kernel v0.96. https://github.com/multipath-tcp/mptcp/blob/mptcp_v0.96/net/mptcp/mptcp_sched.c

[18] MPTCP community. 2024. Multipath TCP Linux Kernel v0.96. https://github.com/multipath-tcp/mptcp/releases

[19] MPTCP GitHub community. 2020. Meta-level OFO queue. https://github.com/multipath-tcp/mptcp/issues/394

[20] MPTCP GitHub community. 2021. Measure MPTCP OFO packets. https://github.com/multipath-tcp/mptcp/issues/460

[21] Mininet GitHub community. 2023. Mininet source code. https://github.com/mininet/mininet

[22] MPTCP GitHub community. 2023. MPTCPv0. https://github.com/multipath-tcp/mptcp

[23] MPTCP GitHub community. 2025. MPTCPv1. https://github.com/multipath-tcp/mptcp_net-next

[24] PyShark GitHub community. 2025. PyShark. https://github.com/KimiNewt/pyshark

[25] Mininet Project Contributors. 2022. *Mininet*. Retrieved September 1, 2025 from https://mininet.org/

[26] Martin Duke. 2023. QUIC Version 2. RFC 9369. doi:10.17487/RFC9369

[27] Martin Duke and Gorry Fairhurst. 2025. Specifying New Congestion Control Algorithms. RFC 9743. doi:10.17487/RFC9743

[28] Wesley Eddy. 2022. Transmission Control Protocol (TCP). RFC 9293. doi:10.17487/RFC9293

[29] Simone Ferlin, Özgü Alay, Olivier Mehani, and Roksana Boreli. 2016. BLEST: Blocking estimation-based MPTCP scheduler for heterogeneous networks. In *2016 IFIP networking conference (IFIP networking) and workshops*. IEEE, 431–439.

[30] S. Ferlin and D. Weber. 2020. BLEST Scheduler Linux Kerner Code. https://github.com/multipath-tcp/mptcp/blob/mptcp_v0.96/net/mptcp/mptcp_blest.c

[31] Sally Floyd. 2000. Congestion Control Principles. RFC 2914. doi:10.17487/RFC2914

[32] Sally Floyd and Mark Allman. 2007. Specifying New Congestion Control Algorithms. RFC 5033. doi:10.17487/RFC5033

[33] Sally Floyd, Hari Balakrishnan, and Mark Allman. 2001. Enhancing TCP's Loss Recovery Using Limited Transmit. RFC 3042. doi:10.17487/RFC3042

[34] Sally Floyd, Mark Handley, and Jitendra Padhye. 2000. A comparison of equation-based and AIMD congestion control. (2000).

[35] Sally Floyd, Jamshid Mahdavi, Matt Mathis, and Dr. Allyn Romanow. 1996. TCP Selective Acknowledgment Options. RFC 2018. doi:10.17487/RFC2018

[36] Alan Ford, Costin Raiciu, Mark J. Handley, Olivier Bonaventure, and Christoph Paasch. 2020. TCP Extensions for Multipath Operation with Multiple Addresses. RFC 8684. doi:10.17487/RFC8684

[37] The Linux Foundation. 2016. TCP Probe. https://wiki.linuxfoundation.org/networking/tcpprobe

[38] The Linux Foundation. 2025. The Linux Kernel Archives. https://www.kernel.org/

[39] Alexander Frommgen, Tobias Erbshäußer, Alejandro Buchmann, Torsten Zimmermann, and Klaus Wehrle. 2016. ReMP TCP: Low latency multipath TCP. In *2016 IEEE international conference on communications (ICC)*. IEEE, 1–7.

[40] Andrei Gurtov, Tom Henderson, Sally Floyd, and Yoshifumi Nishida. 2012. The NewReno Modification to TCP's Fast Recovery Algorithm. RFC 6582. doi:10.17487/RFC6582

[41] S. Ha and S. Hemminger. 2021. CUBIC CCA - Multipath TCP Linux Kernel v0.96. https://github.com/multipath-tcp/mptcp/blob/mptcp_v0.96/net/ipv4/tcp_cubic.c

[42] Sangtae Ha, Injong Rhee, and Lisong Xu. 2008. CUBIC: a new TCP-friendly high-speed TCP variant. *SIGOPS Oper. Syst. Rev.* 42, 5 (July 2008), 64–74. doi:10.1145/1400097.1400105

[43] Jonas Hammar. 2022. *Multipath Packet Scheduling for 5G Systems*. Master's thesis. Department of Mathematics and Computer Science, Karlstad University.

[44] Michio Honda, Yoshifumi Nishida, Lars Eggert, Pasi Sarolahti, and Hideyuki Tokuda. 2009. Multipath congestion control for shared bottleneck. In *Proc. PFLDNeT workshop*, Vol. 357. 378.

[45] J. D. Hunter. 2007. Matplotlib: A 2D graphics environment. *Computing in Science & Engineering* 9, 3 (2007), 90–95. doi:10.1109/MCSE.2007.55

[46] Per Hurtig, Karl-Johan Grinnemo, Anna Brunstrom, Simone Ferlin, Özgü Alay, and Nicolas Kuhn. 2018. Low-latency scheduling in MPTCP. *IEEE/ACM transactions on networking* 27, 1 (2018), 302–315.

[47] Geoff Huston. 2017. *BBR TCP*. Retrieved September 1, 2025 from https://www.potaroo.net/ispcol/2017-05/bbr.pdf

[48] J. Hwang, A. Walid, Q. Peng, and Low S. 2018. BALIA CCA - Multipath TCP Linux Kernel v0.96. https://github.com/multipath-tcp/mptcp/blob/mptcp_v0.96/net/mptcp/mptcp_balia.c

[49] iPerf GitHub Community. 2025. *iPerf: A tool for active measurements of the maximum achievable bandwidth on IP networks*. https://github.com/esnet/iperf

[50] iPerf Team. 2025. iPerf. https://iperf.fr/

[51] Jana Iyengar, Costin Raiciu, Sebastien Barre, Mark J. Handley, and Alan Ford. 2011. Architectural Guidelines for Multipath TCP Development. RFC 6182. doi:10.17487/RFC6182

[52] Jana Iyengar and Martin Thomson. 2021. QUIC: A UDP-Based Multiplexed and Secure Transport. RFC 9000. doi:10.17487/RFC9000

[53] Van Jacobson. 1988. Congestion avoidance and control. *ACM SIGCOMM computer communication review* 18, 4 (1988), 314–329.

[54] V. Jacobson, N. Cardwell, Y. Cheng, and S-H. Yeganeh. 2021. BBR CCA - Multipath TCP Linux Kernel v0.96. https://github.com/multipath-tcp/mptcp/blob/mptcp_v0.96/net/ipv4/tcp_bbr.c

[55] V. Jacobson, N. Cardwell, Y. Cheng, S-H. Yeganeh, and I. Mahmud. 2020. C-MPBBR CCA - Multipath TCP Linux Kernel v0.96. https://github.com/imtiaztee/C-MPBBR/blob/master/mptcp_c-mpbbr.c

[56] Dominik Kaspar. 2011. *Multipath Aggregation of Heterogeneous Access Networks*. Ph. D. Dissertation. Faculty of Mathematics and Natural Sciences, Oslo, Norway.

[57] R. Khalili, N. Gast, and Boudec J-Y. 2021. OLIA CCA - Multipath TCP Linux Kernel v0.96. https://github.com/multipath-tcp/mptcp/blob/mptcp_v0.96/net/mptcp/mptcp_olia.c

[58] Ramin Khalili, Nicolas Gast, Miroslav Popovic, and Jean-Yves Le Boudec. 2014. *Opportunistic Linked-Increases Congestion Control Algorithm for MPTCP*. Internet-Draft draft-khalili-mptcp-congestion-control-05. Internet Engineering Task Force. https://datatracker.ietf.org/doc/draft-khalili-mptcp-congestion-control/05/ Work in Progress.

[59] Ramin Khalili, Nicolas Gast, Miroslav Popovic, and Jean-Yves Le Boudec. 2013. MPTCP is not pareto-optimal: performance issues and a possible solution. *IEEE/ACM Trans. Netw.* 21, 5 (Oct. 2013), 1651–1665. doi:10.1109/TNET.2013.2274462

[60] Bruno YL Kimura, Demetrius CSF Lima, and Antonio AF Loureiro. 2020. Packet scheduling in multipath TCP: Fundamentals, lessons, and opportunities. *IEEE Systems Journal* 15, 1 (2020), 1445–1457.

[61] Nicolas Kuhn, Emmanuel Lochin, Ahlem Mifdaoui, Golam Sarwar, Olivier Mehani, and Roksana Boreli. 2014. DAPS: Intelligent delay-aware packet scheduling for multipath transport. In *2014 IEEE international conference on communications (ICC)*. IEEE, 1222–1227.

[62] Bob Lantz, Brandon Heller, and Nick McKeown. 2010. A network in a laptop: rapid prototyping for software-defined networks. In *Proceedings of the 9th ACM SIGCOMM Workshop on Hot Topics in Networks* (Monterey, California) *(Hotnets-IX)*. Association for Computing Machinery, New York, NY, USA, Article 19, 6 pages. doi:10.1145/1868447.1868466

[63] Y. Lim and D. Weber. 2020. ECF Scheduler Linux Kerner Code. https://github.com/multipath-tcp/mptcp/blob/mptcp_v0.96/net/mptcp/mptcp_ecf.c

[64] Yeon-sup Lim, Erich M Nahum, Don Towsley, and Richard J Gibbens. 2017. ECF: An MPTCP path scheduler to manage heterogeneous paths. In *Proceedings of the 13th international conference on emerging networking experiments and technologies*. 147–159.

[65] Yanmei Liu, Yunfei Ma, Quentin De Coninck, Olivier Bonaventure, Christian Huitema, and Mirja Kühlewind. 2025. *Multipath Extension for QUIC*. Internet-Draft draft-ietf-quic-multipath-16. Internet Engineering Task Force. https://datatracker.ietf.org/doc/draft-ietf-quic-multipath/16/ Work in Progress.

[66] T. Lubna and I. Mahmud. 2022. Delay-Data Rate Scheduler Linux Kerner Code. https://github.com/imtiaztee/DRS

[67] Tabassum Lubna, Imtiaz Mahmud, and You-Ze Cho. 2022. Low latency and high data rate (LLHD) scheduler: a multipath TCP scheduler for dynamic and heterogeneous networks. *Sensors* 22, 24 (2022), 9869.

[68] Imtiaz Mahmud, Tabassum Lubna, Yeong-Jun Song, and You-Ze Cho. 2020. Coupled multipath BBR (C-MPBBR): A efficient congestion control algorithm for multipath TCP. *IEEE access* 8 (2020), 165497–165511.

[69] Christoph Paasch. 2019. RR - Multipath TCP Linux Kernel v0.96. https://github.com/multipath-tcp/mptcp/blob/mptcp_v0.96/net/mptcp/mptcp_rr.c

[70] C. Paasch and S. Barré. 2018. LIA CCA - Multipath TCP Linux Kernel v0.96. https://github.com/multipath-tcp/mptcp/blob/mptcp_v0.96/net/mptcp/mptcp_coupled.c

[71] Christoph Paasch, Simone Ferlin, Ozgu Alay, and Olivier Bonaventure. 2014. Experimental evaluation of multipath TCP schedulers. In *Proceedings of the 2014 ACM SIGCOMM workshop on Capacity sharing workshop*. ACM, New York, USA, 27–32.

[72] Christoph Paasch, Ramin Khalili, and Olivier Bonaventure. 2013. On the benefits of applying experimental design to improve multipath TCP. In *Proceedings of the Ninth ACM Conference on Emerging Networking Experiments and Technologies* (Santa Barbara, California, USA) *(CoNEXT '13)*. Association for Computing Machinery, New York, NY, USA, 393–398. doi:10.1145/2535372.2535403

[73] Dr. Vern Paxson and Mark Allman. 2000. Computing TCP's Retransmission Timer. RFC 2988. doi:10.17487/RFC2988

[74] Qiuyu Peng, Anwar Walid, Jaehyun Hwang, and Steven H Low. 2014. Multipath TCP: Analysis, design, and implementation. *IEEE/ACM Transactions on networking* 24, 1 (2014), 596–609.

[75] René Pfeiffer. 2007. Measuring TCP Congestion Windows. https://linuxgazette.net/136/pfeiffer.html

[76]  Michele Polese, Federico Chiariotti, Elia Bonetto, Filippo Rigotto, Andrea Zanella, and Michele Zorzi. 2019. A survey on recent advances in transport
      layer protocols. *IEEE Communications Surveys & Tutorials* 21, 4 (2019), 3584–3608.

[77]  Python Core Team. 2024. *Python: A dynamic, open source programming language.* Python Software Foundation. https://www.python.org/ Python
      version 3.10.

[78]  Costin Raiciu, Mark J. Handley, and Damon Wischik. 2011. Coupled Congestion Control for Multipath Transport Protocols. RFC 6356. doi:10.17487/
      RFC6356

[79]  Costin Raiciu, Christoph Paasch, Sebastien Barre, Alan Ford, Michio Honda, Fabien Duchene, Olivier Bonaventure, and Mark Handley. 2012.
      How hard can it be? designing and implementing a deployable multipath {TCP}. In *9th USENIX symposium on networked systems design and
      implementation (NSDI 12)*. USENIX, 399–412.

[80]  Swetank Kumar Saha, Shivang Aggarwal, Rohan Pathak, Dimitrios Koutsonikolas, and Joerg Widmer. 2019. MuSher: An agile multipath-TCP
      scheduler for dual-band 802.11 ad/AC wireless LANs. In *The 25th Annual International Conference on Mobile Computing and Networking*. 1–16.

[81]  Matt Sargent, Jerry Chu, Dr. Vern Paxson, and Mark Allman. 2011. Computing TCP's Retransmission Timer. RFC 6298. doi:10.17487/RFC6298

[82]  Stefan Savage, Neal Cardwell, David Wetherall, and Tom Anderson. 1999. TCP congestion control with a misbehaving receiver. *ACM SIGCOMM
      Computer Communication Review* 29, 5 (1999), 71–78.

[83]  Michael Scharf and Alan Ford. 2013. Multipath TCP (MPTCP) Application Interface Considerations. RFC 6897. doi:10.17487/RFC6897

[84]  Reinhard Schrage, Gilles Forget, Ruediger Geib, and Barry Constantine. 2011. Framework for TCP Throughput Testing. RFC 6349. doi:10.17487/
      RFC6349

[85]  Anwar Walid, Qiuyu Peng, Jaehyun Hwang, and Steven H. Low. 2016. *Balanced Linked Adaptation Congestion Control Algorithm for MPTCP.*
      Internet-Draft draft-walid-mptcp-congestion-control-04. Internet Engineering Task Force. https://datatracker.ietf.org/doc/draft-walid-mptcp-
      congestion-control/04/ Work in Progress.

[86]  Wikipedia contributors. 2025. TCP congestion control — Wikipedia, The Free Encyclopedia. https://en.wikipedia.org/w/index.php?title=TCP_
      congestion_control&oldid=1308091666 [Online; accessed 28-August-2025].

[87]  Damon Wischik, Mark Handley, and Marcelo Bagnulo Braun. 2008. The resource pooling principle. *ACM SIGCOMM Computer Communication
      Review* 38, 5 (2008), 47–52.

[88]  Damon Wischik, Costin Raiciu, Adam Greenhalgh, and Mark Handley. 2011. Design, implementation and evaluation of congestion control for
      multipath {TCP}. In *8th USENIX Symposium on Networked Systems Design and Implementation (NSDI 11)*.

[89]  Lisong Xu, Sangtae Ha, Injong Rhee, Vidhi Goel, and Lars Eggert. 2023. CUBIC for Fast and Long-Distance Networks. RFC 9438. doi:10.17487/RFC9438

[90]  Mingwei Xu, Yu Cao, and Enhuan Dong. 2017. *Delay-based Congestion Control for MPTCP.* Internet-Draft draft-xu-mptcp-congestion-control-05.
      Internet Engineering Task Force. https://datatracker.ietf.org/doc/draft-xu-mptcp-congestion-control/05/ Work in Progress.

[91]  Fan Yang, Qi Wang, and Paul D Amer. 2014. Out-of-order transmission for in-order arrival scheduling for multipath TCP. In *2014 28th international
      conference on advanced information networking and applications workshops*. IEEE, 749–752.

[92]  Ming Zhang, Junwen Lai, Arvind Krishnamurthy, Larry L Peterson, and Randolph Y Wang. 2004. A Transport Layer Approach for Improving
      End-to-End Performance and Robustness Using Redundant Paths.. In *USENIX Annual Technical Conference, General Track*. 99–112.

## A  Packet Scheduling Algorithms

---

**Algorithm 1:** minRTT

---

**Input:** Set of available subflows: $\mathcal{S}$,  set of paths: $\mathcal{R}$,  $sf_i \in \mathcal{S}$, $i \in \mathcal{R}$

**Output:** Best subflow, $sf_{best} \in \mathcal{S}$

1  $srtt_{min} \leftarrow 0xffffffff$ ;

2  $sf_{best} \leftarrow NULL$ ;

3  **for** *each* $sf_i \in \mathcal{S}$ **do**

4      **if** $srtt_{sf_i} < srtt_{min}$ **then**

5          $srtt_{min} \leftarrow srtt_{sf_i}$ ;

6          $sf_{best} \leftarrow sf_i$ ;

7      **end if**

8  **end for**

9  **return** $sf_{best}$

---

The theoretical aspects of *minRTT* packet scheduling algorithm are provided in 2.2.1.

---

**Algorithm 2:** BLEST

---

**Input:** Set of available subflows: $\mathcal{S}$,  set of paths: $\mathcal{R}$,  $sf_i \in \mathcal{S}$, $i \in \mathcal{R}$,  $srtt_f < srtt_s$

**Output:** Best subflow, $sf_{best} \in \mathcal{S}$

1  $x_f \leftarrow minRTT()$ ;

2  **if** *fastest subflow* $x_f$ *is available* **then**

3      $sf_{best} \leftarrow x_f$ ;

4  **else if** *slower subflow* $x_s$ *is available* **then**

5      $rtt_s = srtt_s/srtt_f$ ;

6      $X = MSS_f * (CWND_f + (rtt_s - 1)/2) * rtt_s$ ;

7      **if** $X * \lambda \leq MPTCP_{SW} - MSS_s * (inflight_s + 1)$ **then**

8          $sf_{best} \leftarrow x_s$ ;

9      **else**

10          $sf_{best} \leftarrow NULL$ ;

11      **end if**

12  **else**

13      $sf_{best} \leftarrow NULL$ ;                                   `/* no available subflow */`

14  **end if**

15  **return** $sf_{best}$

---

The theoretical aspects of *BLEST* packet scheduling algorithm are provided in 2.2.2.

---

**Algorithm 3:** ECF

---

**Input:** Set of available subflows: $\mathcal{S}$, set of paths: $\mathcal{R}$, $sf_i \in \mathcal{S}$, $i \in \mathcal{R}$, $srtt_f < srtt_s$

**Output:** Best subflow, $sf_{best} \in \mathcal{S}$

1   $x_f \leftarrow minRTT()$ ;                       `/* find the fastest subflow xf using minRTT algorithm */`

2   **if** *fastest subflow* $x_f$ *is available* **then**

3      $sf_{best} \leftarrow x_f$ ;

4   **else**

5      $x_s \leftarrow minRTT()$ ;            `/* find the second fastest subflow xs using minRTT algorithm */`

6      $n = 1 + \frac{k}{CWND_f}$ ;

7      $\delta = max(\sigma_f, \sigma_s)$ ;

8      **if** $n * RTT_f < (1 + waiting * \beta)(RTT_s + \delta)$ **then**

9          **if** $\frac{k}{CWND_s} * RTT_s \geq 2 * RTT_f + \delta$ **then**

10             $waiting = 1$ ;                                   `/* Wait for xf */`

11             $sf_{best} \leftarrow NULL$ ;                        `/* no available subflow */`

12          **else**

13             $sf_{best} \leftarrow x_s$ ;

14          **end if**

15      **else**

16          $waiting = 0$ ;

17          $sf_{best} \leftarrow x_s$ ;

18      **end if**

19   **end if**

20   **return** $sf_{best}$

---

The theoretical aspects of *ECF* packet scheduling algorithm are provided in 2.2.3.

---

**Algorithm 4:** RR (num_segments $\geq$ 1, cwnd_limited=true)

---

**Input:** Set of available subflows: $\mathcal{S}$, set of paths: $\mathcal{R}$, $sf_i \in \mathcal{S}$, $i \in \mathcal{R}$

1   Get *num_segments* from MPTCP send buffer ;

2   **for** *each* $sf_i \in \mathcal{S}$ **do**

3      **if** $sf_{i\_}cwnd - sf_{i\_}inflight \geq num\_segments * MSS_i$ **then**

4          send (*num_segments* * $MSS_i$) bytes on $sf_i$ ;           `/* sfi has sufficient cwnd space */`

5      **else**

6          continue ;                       `/* skip sfi in case of insufficient cwnd space */`

7      **end if**

8   **end for**

---

The theoretical aspects of *RR* packet scheduling algorithm are provided in 2.2.4.

---

**Algorithm 5:** LLHD

---

**Input:** Set of subflows: $\mathcal{S}$, set of paths: $\mathcal{R}$, $sf_i \in \mathcal{S}$, $i \in \mathcal{R}$

**Output:** Best subflow, $sf_{best} \in \mathcal{S}$

**Initialization:**

$sf_{best} = NULL$ ;

$\gamma_{max} = 0$ ;

$\beta = 0.001$ ;               /* $\beta = 0.001$ in paper [67]; however, actual value in code is $\beta = 10$ [66] */

$RTT_{max} = 9999999$ ;

$GP_{max} = 9999$ ;

Upon reception of ACK:

1  **for** *each* $sf_i \in \mathcal{S}$ **do**

2      **if** $sf_i$ *is backup* **then**

3          $sf_i = NULL$ ;

4          $sf_{best} \leftarrow sf_i$ ;

5      **end if**

6      **if** $sf_i$ *is unavailable* **then**

7          $sf_i = NULL$ ;

8          $sf_{best} \leftarrow sf_i$ ;

9      **end if**

10     **if** $sf_i$ *is temp_unavailable* **then**

11         $sf_i = NULL$ ;

12         $sf_{best} \leftarrow sf_i$ ;

13     **end if**

14     $\gamma\_curr = (GP\_sf_i/GP_{max}) + \beta \times (RTT_{max}/RTT\_sf_i)$ ;

15     **if** $sf_i\_CWND\_available$ **and** $\gamma\_curr > \gamma\_max$ **then**

16         $\gamma_{max} \leftarrow \gamma\_curr$ ;

17         $sf_{best} \leftarrow sf_i$ ;

18     **end if**

19 **end for**

20 **return** $sf_{best}$

---

The theoretical aspects of *LLHD* packet scheduling algorithm are provided in 2.2.5.

---

**Algorithm 6:** ReMP TCP (redundant)

---

**Input:** Set of subflows: $\mathcal{S}$, set of paths: $\mathcal{R}$, $sf_i \in \mathcal{S}$, $i \in \mathcal{R}$

**Output:** $sf_i$, $sf_j \in \mathcal{S}$

1 **for** *each segment to be transmitted* **do**
2      **for** *each $sf_k \in \mathcal{S}$* **do**
3          **if** *$sf_k$ is available* **then**
4              $sf_i \leftarrow sf_k$ ;
5              **return** $sf_i$ ;
6          **end if**
7      **end for**
8      **for** *each $sf_k \in \mathcal{S}$* **do**
9          **if** *$sf_k \neq sf_i$ and $sf_k$ is available* **then**
10              $sf_j \leftarrow sf_k$ ;
11              **return** $sf_j$ ;
12          **end if**
13      **end for**
14 **end for**
15 redundant**return** NULL ;                     /* either one or both subflows unavailable */

---

The theoretical aspects of *ReMP TCP* packet scheduling algorithm are provided in 2.2.6.

## B    Congestion Control Algorithms

---

**Algorithm 7:** CUBIC (v2.2)

---

Initialization:
     $tcp\_friendliness \longleftarrow 1$, $\beta \longleftarrow 0.2$
     $fast\_convergence \longleftarrow 1$, $C \longleftarrow 0.4$
     $cubic\_reset()$
On each ACK:
**begin**
1      **if** *dMin* **then** $dMin \longleftarrow min(dMin, RTT)$
2      **else** $dMin \longleftarrow RTT$
3      **if** $cwnd \leq ssthresh$ **then** $cwnd \longleftarrow cwnd + 1$
4      **else**
5          $cnt \longleftarrow cubic\_update()$
6          **if** $cwnd\_cnt > cnt$ **then**
7              $cwnd \longleftarrow cwnd + 1$,   $cwnd\_cnt \longleftarrow 0$
8          **else** $cwnd\_cnt \longleftarrow cwnd\_cnt + 1$

---

Packet loss:

**begin**

9     $epoch\_start \longleftarrow 0$

10     **if** $cwnd < W_{last\_max}$ **and** $fast\_convergence$ **then**

11       $W_{last\_max} \longleftarrow cwnd * \frac{(2-\beta)}{2}$

12     **else** $W_{last\_max} \longleftarrow cwnd$

13     $ssthresh \longleftarrow cwnd \longleftarrow cwnd * (1 - \beta)$

Timeout:

**begin**

14     $cubic\_reset()$

15 $cubic\_update()$ :

**begin**

16     $ack\_cnt \longleftarrow ack\_cnt + 1$

17     **if** $epoch\_start \leq 0$ **then**

18       $epoch\_start \longleftarrow tcp\_time\_stamp$

19       **if** $cwnd < W_{last\_max}$ **then**

20         $K \longleftarrow \sqrt[3]{\frac{W_{last\_max} - cwnd}{C}}$

21         $origin\_point \longleftarrow W_{last\_max}$

22       **else**

23         $K \longleftarrow 0$

24         $origin\_point \longleftarrow cwnd$

25       $ack\_cnt \longleftarrow 1$

26       $W_{tcp} \longleftarrow cwnd$

27     $t \longleftarrow tcp\_time\_stamp + dMin - epoch\_start$

28     $target \longleftarrow origin\_point + C(t - K)^3$

29     **if** $target > cwnd$ **then** $cnt \longleftarrow \frac{cwnd}{target - cwnd}$

30     **else** $cnt \longleftarrow 100 * cwnd$

31     **if** $tcp\_friendliness$ **then** $cubic\_tcp\_friendliness()$

32 $cubic\_tcp\_friendliness()$ :

**begin**

33     $W_{tcp} \longleftarrow W_{tcp} + \frac{3\beta}{2-\beta} * \frac{ack\_cnt}{cwnd}$

34     $ack\_cnt \longleftarrow 0$

35     **if** $W_{tcp} > cwnd$ **then**

36       $max\_cnt \longleftarrow \frac{cwnd}{W_{tcp} - cwnd}$

37       **if** $cnt > max\_cnt$ **then** $cnt \longleftarrow max\_cnt$

38 $cubic\_reset()$ :

**begin**

39     $W_{last\_max} \longleftarrow 0,\ epoch\_start \longleftarrow 0,\ origin\_point \longleftarrow 0$

40     $dMin \longleftarrow 0,\ W_{tcp} \longleftarrow 0,\ K \longleftarrow 0,\ ack\_cnt \longleftarrow 0$

The theoretical aspects of *Cubic* congestion control algorithm are provided in 2.3.1.

---

**Algorithm 8:** Coupled (LIA) - Congestion Avoidance phase (AIMD)

---

**Input:** Set of subflows: $\mathcal{S}$, set of paths: $\mathcal{R}$, $sf_i \in \mathcal{S}$, $i \in \mathcal{R}$, $cwnd_i$, $cwnd_{total}$, $RTT_i$, $MSS_i$

**Output:** $cwnd_i$

1 **for** $sf_i \in \mathcal{S}$ **do**

    **On each ACK received on** $sf_i$:

2       $\alpha \longleftarrow cwnd_{total} * \frac{max(cwnd_i/RTT_i^2)}{(\sum_i (cwnd_i/RTT_i))^2}$

3       $cwnd_i \longleftarrow cwnd_i + min(\frac{\alpha * bytes\_acked * MSS_i}{cwnd_{total}}, \frac{bytes\_acked * MSS_i}{cwnd_i})$

                                                   `/* in the above min(A, B) computation, A:`
                                                      `denotes the computed increase of the`
                                                     `multipath subflow, B: denotes the increase`
                                                     `TCP would get on the same path i */`

    **On packet loss event on** $sf_i$:

4       $cwnd_i \longleftarrow \frac{cwnd_i}{2}$

5       $fast\_retransmit()$;                           `/* legacy TCP NewReno */`

6       $fast\_recovery()$;

7 **return** $cwnd_i$

---

The theoretical aspects of *Coupled (LIA)* congestion control algorithm are provided in 2.3.2.

---

**Algorithm 9:** Opportunistic Linked-Increases Algorithm (OLIA) - Congestion Avoidance phase (AIMD)

---

**Input:**

    set of paths available to user u: $R_u$, $r \in R_u$     `/*` $R_u$`: the set of` **`all_paths`**`,` $|R_u|$`: the number`
                                                      `of paths available to user` $u$ `at time` $t$ `*/`

    set of paths with the largest window sizes: $\mathcal{M}$,     `/*` $\mathcal{M}$`: the set of` **`max_w_paths`** `of user` $u$ `at`
                                                       `time` $t$`,` $|\mathcal{M}|$`: the number of paths in` $\mathcal{M}$ `*/`

    set of presumably best paths: $\mathcal{B}$,     `/*` $\mathcal{B}$`: the set of presumably` **`best_paths`** `for`
                                                       `user` $u$ `at time` $t$`,` $|\mathcal{B}|$`: the number of paths`
                                                       `in` $\mathcal{B}$ `*/`

    set of collected paths: $\mathcal{B}\backslash\mathcal{M}$,     `/*` $\mathcal{B}\backslash\mathcal{M}$`: the set of` **`collected_paths`** `for user`
                                                       `u` `at time` $t$ `- all paths belonging to`
                                                       `best_paths, but not among those with`
                                                      `largest` `cwnd` `*/`

    congestion window of path $r$: $w_r$

    round-trip time of path $r$: $rtt_r$

    a subflow established over path $r$: $sf_r$

    set of subflows established: $S$, where $sf_r \in S$

---

**Output:** $w_r$

1 **for** $sf_r \in S$ **do**

**On each ACK received on** $sf_r$:

2
$$a_r = \begin{cases} \frac{1/|R_u|}{|\mathcal{B}\backslash\mathcal{M}|}, & \text{if } r \in \mathcal{B}\backslash\mathcal{M} \neq \emptyset \\ -\frac{1/|R_u|}{|\mathcal{M}|}, & \text{if } r \in \mathcal{M} \text{ and } \mathcal{B}\backslash\mathcal{M} \neq \emptyset \\ 0, & \text{otherwise.} \end{cases}$$
```
/* calculate the increase parameter aᵣ //
```
also note that $\sum_{r \in R_u} a_r = 0$   `*/`

3
$$w_r \longleftarrow w_r + \left( \frac{w_r/rtt_r^2}{\left(\sum_{p\in R_u} w_p/rtt_p\right)^2} + \frac{a_r}{w_r} \right)$$
```
/* window increase measured in packets //
   multiply by (MSSᵣ * bytes_acked) to measure
   in bytes  */
```

**On packet loss event on** $sf_r$:

4     $w_r \longleftarrow \frac{w_r}{2}$

5     $fast\_retransmit();$                                 `/* legacy TCP NewReno  */`

6     $fast\_recovery();$

7 **return** $w_r$

---

The theoretical aspects of *OLIA* congestion control algorithm are provided in 2.3.3.

---

**Algorithm 10:** Balanced Linked Adaptation (BALIA) - Congestion Avoidance phase (AIMD)

**Input:**

set of available paths: $R, \ r \in R$

congestion window of path $r$: $w_r$

round-trip time of path $r$: $rtt_r$

a subflow established over path $r$: $sf_r$

set of subflows established: $S$, where $sf_r \in S$

**Output:** $w_r$

1 **for** $sf_r \in S$ **do**

**On each ACK received on** $sf_r$:

2     $a_r := \frac{\max\{x_k\}}{x_r}$
```
/* calculate the increase parameter aᵣ //
   xᵣ := wᵣ/rttᵣ // also note that in case of
   single path aᵣ = 1  */
```

3     $w_r \longleftarrow w_r + \frac{x_r}{rtt_r(\sum x_k)^2} \left(\frac{1+a_r}{2}\right)\left(\frac{4+a_r}{5}\right)$

**On packet loss event on** $sf_r$:

4     $w_r \longleftarrow w_r - \frac{w_r}{2}\min\{a_r, 1.5\}$

5     $fast\_retransmit();$                                 `/* legacy TCP NewReno  */`

6     $fast\_recovery();$

7 **return** $w_r$

---

The theoretical aspects of *BALIA* congestion control algorithm are provided in 2.3.4.

---

**Algorithm 11:** Weighted Vegas (wVegas)

---

1 **Initialization:**

3   $total\_alpha \leftarrow 10$ ;    /* namely $\alpha_s$: preconfigured parameter coupling subflows of flow s */

4   **for** $r \in R_s$ **do**

5    $alpha[r] \leftarrow 2$ ;

6    $equilibrium\_rates[r] \leftarrow 0$ ;

7    $queue\_delays[r] \leftarrow 0$ ;

 

**On the end of round for subflow** $r$**:**

  /* average RTT estimated in the last round - used instead of smoothed RTT for faster

   reaction to congestion                   */

8   $rtt \leftarrow sampled\_rtts[r]/sampled\_num[r]$ ;

9   $diff \leftarrow cwnd[r] \times (rtt - baseRTT[r])/rtt$ ;

  /* tweak weights and alphas                         */

10   **if** $diff \geq alpha[r]$ **then**

11    $equilibrium\_rates[r] \leftarrow cwnd[r]/rtt$ ;

12    Adjust_Weights() ;

13    $alpha[r] \leftarrow weights[r] \times total\_alpha$ ;

14    $alpha[r] \leftarrow \max\{2, alpha[r]\}$ ; // lower bound

  /* window adjustment                            */

15   **if** $diff < alpha[r]$ **then**

16    $cwnd[r] \leftarrow cwnd[r] + 1$ ;

17   **else if** $diff > alpha[r]$ **then**

18    $cwnd[r] \leftarrow cwnd[r] - 1$ ;

  /* try to drain link queues if needed                */

19   $q \leftarrow rtt - baseRTT[r]$ ; // current queuing delay

20   **if** $queue\_delays[r] = 0$ **or** $queue\_delays[r] > q$ **then**

21    $queue\_delays[r] \leftarrow q$ ;

22   **if** $q \geq 2 \times queue\_delays[r]$ **then**    /* cwnd backoff once queuing delay exceeds threshold */

23    $backoff\_factor \leftarrow 0.5 \times baseRTT[r]/rtt$

24    $cwnd[r] \leftarrow cwnd[r] \times backoff\_factor$

25    $queue\_delays[r] \leftarrow 0$

26   $cwnd[r] \leftarrow \max\{2, cwnd[r]\}$ // lower bound

 

**Adjust_Weights():**

27   $total\_rate \leftarrow \sum equilibrium\_rates$ ;

28   **for** $r \in R_s$ **do**

29    **if** $equilibrium\_rates[r] \neq 0$ **then**

30     $weights[r] \leftarrow equilibrium\_rates[r]/total\_rate$ ;

 

**On packet loss for subflow** $r$**:**

31   $equilibrium\_rates[r] \leftarrow 0$ ;

32   $queue\_delays[r] \leftarrow 0$ ;

---

The theoretical aspects of *wVegas* congestion control algorithm are provided in 2.3.5.

---

**Algorithm 12:** Bottleneck Bandwidth and Round-trip propagation time (BBR) - [theoretical part in 2.3.6]

---

**1  Initialization:**
2     BBR_On_Connection_Init():
3        BBR_Init();

  **On ACK arrival:**
4     BBR_Update_On_ACK():
5        *BBR_Update_Model_And_State*() ;
6        *BBR_Update_Control_Parameters*() ;

7     BBR_Update_Model_And_State():
8        *BBR_Update_BtlBw*() ;
9        *BBR_Check_Cycle_Phase*() ;
10       *BBR_Check_Full_Pipe*() ;
11       *BBR_Check_Drain*() ;
12       *BBR_Update_RTprop*() ;
13       *BBR_Check_Probe_RTT*() ;

14    BBR_Update_Control_Parameters():
15       *BBR_Set_Pacing_Rate*() ;
16       *BBR_Set_Send_Quantum*() ;
17       *BBR_Set_Cwnd*() ;

  **On packet transmission:**
18    BBR_On_Transmit():
19       *BBR_Handle_Restart_From_Idle*() ;

  **Initialization function definitions:**
20    BBR_Init():
21       *init_windowed_max_filter*($filter = BBR.BtlBwFilter$, $value = 0$, $time = 0$) ;
22       $BBR.rtprop = SRTT\ ?\ SRTT\ :\ Inf$ ;
23       $BBR.rtprop\_stamp = Now()$ ;
24       $BBR.probe\_rtt\_done\_stamp = 0$ ;
25       $BBR.probe\_rtt\_round\_done = false$ ;
26       $BBR.packet\_conservation = false$ ;
27       $BBR.prior\_cwnd = 0$ ;
28       $BBR.idle\_restart = false$ ;
29       *BBR_Init_Round_Counting*() ;
30       *BBR_Init_Full_Pipe*() ;
31       *BBR_Init_Pacing_Rate*() ;
32       *BBR_Enter_Startup*() ;
33    BBR_Init_Round_Counting():
34       $BBR.next\_round\_delivered = 0$ ;
35       $BBR.round\_start = false$ ;
36       $BBR.round\_count = 0$ ;

---

**Initialization function definitions (cont.):**

37    BBR_Init_Full_Pipe():

38      $BBR.filled\_pipe = false$ ;

39      $BBR.full\_bw = 0$ ;

40      $BBR.full\_bw\_count = 0$ ;

41    BBR_Init_Pacing_Rate():

42      $nominal\_bandwidth = InitialCwnd/(SRTT \; ? \; SRTT \; : \; 1ms)$ ;

43      $BBR.pacing\_rate = BBR.pacing\_gain * nominal\_bandwidth$ ;

44    BBR_Enter_Startup():

45      $BBR.state = Startup$ ;

46      $BBR.pacing\_gain = BBRHighGain$ ;

47      $BBR.cwnd\_gain = BBRHighGain$ ;

**ACK-arrival function definitions (1/4):**

48    BBR_Update_BtlBw():

49      $BBR\_Update\_Round()$ ;

50      **if** $rs.delivery\_rate \geq BBR.BtlBw$ **or not** $rs.is\_app\_limited$ **then**

51        $BBR.BtlBw = update\_windowed\_max\_filter(filter = BBR.BtlBwFilter, value = rs.delivery\_rate,$
         $time = BBR.round\_count, window\_length = BtlBwFilterLen)$ ;

52    BBR_Update_Round():

53      $BBR.delivered+ = packet.size$ ;

54      **if** $packet.delivered \geq BBR.next\_round\_delivered$ **then**

55        $BBR.next\_round\_delivered = BBR.delivered$ ;

56        $BBR.round\_count + +$ ;

57        $BBR.round\_start = true$ ;

58      **else**

59        $BBR.round\_start = false$ ;

60    BBR_Check_Cycle_Phase():

61      **if** $BBR.sate == ProbeBW$ **and** $BBR\_Is\_Next\_Cycle\_Phase()$ **then**

62        $BBR\_Advance\_Cycle\_Phase()$ ;

63    BBR_Is_Next_Cycle_Phase():

64      $is\_full\_length = (Now() - BBR.cycle\_stamp) > BBR.RTprop$;

65      **if** $BBR.pacing\_gain == 1$ **then**

66        **return** $is\_full\_length$ ;

67      **if** $BBR.pacing\_gain > 1$ **then**

68        **return** $is\_full\_length$ **and** $(packets\_lost > 0$ **or**
         $prior\_inflight \geq BBR\_Inflight(BBR.pacing\_gain))$ ;

69      **else**                                      /* (BBR.pacing_gain < 1) */

70        **return** $is\_full\_length$ **or** $prior\_inflight \leq BBR\_Inflight(1)$ ;

**ACK-arrival function definitions (cont. 2/4):**

```
71   BBR_Advance_Cycle_Phase():
72       BBR.cycle_stamp = Now() ;
73       BBR.cycle_index = (BBR.cycle_index + 1) % BBRGainCycleLen;
74       pacing_gain_cycle = [5/4, 3/4, 1, 1, 1, 1, 1, 1] ;
75       BBR.pacing_gain = pacing_gain_cycle[BBR.cycle_index] ;
76   BBR_Check_Full_Pipe():
77       if BBR.filled_pipe or not BBR.round_start or rs.is_app_limited then
78           return ;                                /* no need to check for a full pipe now */
79       if BBR.BtlBw ≥ BBR.full_bw * 1.25 then              /* BBR.BtlBw still growing? */
80           BBR.full_bw = BBR.BtlBw ;                       /* record new baseline level */
81           BBR.full_bw_count = 0 ;
82           return ;
83       BBR.full_bw_count + + ;                             /* another round w/o much growth */
84       if BBR.full_bw_count ≥ 3 then
85           BBR.filled_pipe = true ;
86   BBR_Check_Drain():
87       if BBR.state == Startup and BBR.filled_pipe then
88           BBR_Enter_Drain() ;
89       if BBR.state == Drain and packets_in_flight ≤ BBR_Inflight(1.0) then
90           BBR_Enter_ProbeBW() ;                           /* we estimate queue is drained */
91   BBR_Update_RTprop():
92       BBR.rtprop_expired = Now() > BBR.rtprop_stamp + RTpropFilterLen ;
93       if packet.rtt ≥ 0 and (packet.rtt ≤ BBR.RTprop or BBR.rtprop_expired) then
94           BBR.RTprop = packet.rtt ;
95           BBR.rtprop_stamp = Now() ;
96   BBR_Check_Probe_RTT():
97       if BBR.state ! = ProbeRTT and BBR.rtprop_expired and not BBR.idle_restart then
98           BBR_Enter_Probe_RTT() ;
99           BBR_Save_Cwnd() ;
100          BBR.probe_rtt_done_stamp = 0 ;
101          if BBR.state == ProbeRTT then
102              BBR_Handle_Probe_RTT() ;
103          BBR.idle_restart = false ;
104  BBR_Enter_Probe_RTT():
105      BBR.state = ProbeRTT ;
106      BBR.pacing_gain = 1 ;
107      BBR.cwnd_gain = 1 ;
```

**ACK-arrival function definitions (cont. 3/4):**

108    BBR_Handle_Probe_RTT():

     `/* Ignore low rate samples during ProbeRTT:`                        `*/`

109      $C.app\_limited = (BW.delivered + packets\_in\_flight)\,?\,:\,1$ ;

110      **if** $BBR.probe\_rtt\_done\_stamp == 0$ **and** $packets\_in\_flight \leq BBRMinPipeCwnd$ **then**

111        $BBR.probe\_rtt\_done\_stamp = Now() + ProbeRTTDuration$ ;

112        $BBR.probe\_rtt\_round\_done = false$ ;

113        $BBR.next\_round\_delivered = BBR.delivered$ ;

114      **else if** $BBR.probe\_rtt\_done\_stamp\,! = 0$ **then**

115        **if** $BBR.round\_start$ **then**

116          $BBR.probe\_rtt\_round\_done = true$ ;

117        **if** $BBR.probe\_rtt\_round\_done$ **and** $Now() > BBR.probe\_rtt\_done\_stamp$ **then**

118          $BBR.rtprop\_stamp = Now()$ ;

119          $BBR\_Restore\_Cwnd()$ ;

120          $BBR\_Exit\_Probe\_RTT()$ ;

121    BBR_Exit_Probe_RTT():

122      **if** $BBR.filled\_pipe$ **then**

123        $BBR\_Enter\_Probe\_BW()$ ;

124      **else**

125        $BBR\_Enter\_Startup()$ ;

126    BBR_Enter_Probe_BW():

127      $BBR.state = ProbeBW$ ;

128      $BBR.pacing\_gain = 1$ ;

129      $BBR.cwnd\_gain = 2$ ;

130      $BBR.cycle\_index = BBRGainCycleLen - 1 - random\_int\_in\_range(0..6)$ ;

131      $BBR\_Advance\_Cycle\_Phase()$ ;

132    BBR_Set_Pacing_Rate():

133      $BBR\_Set\_Pacing\_Rate\_With\_Gain(BBR.pacing\_gain)$ ;

134    BBR_Set_Pacing_Rate_With_Gain(pacing_gain):

135      $rate = pacing\_gain * BBR.BtlBw$ ;

136      **if** $BBR.filled\_pipe$ **or** $rate > BBR.pacing\_rate$ **then**

137        $BBR.pacing\_rate = rate$ ;

138    BBR_Set_Send_Quantum():

139      **if** $BBR.pacing\_rate < 1.2\,Mbps$ **then**

140        $BBR.send\_quantum = 1 * MSS$ ;

141      **else if** $BBR.pacing\_rate < 24\,Mbps$ **then**

142        $BBR.send\_quantum = 2 * MSS$ ;

143      **else**

144        $BBR.send\_quantum = min(BBR.pacing\_rate * 1ms, 64KBytes)$ ;

**ACK-arrival function definitions (cont. 4/4):**

145    BBR_Set_Cwnd():
146        *BBR_Update_Target_Cwnd()* ;
147        *BBR_Modulate_Cwnd_For_Recovery()* ;
148        **if not** *BBR.packet_conservation* **then**
149            **if** *BBR.filled_pipe* **then**
150                $cwnd = min(cwnd + packets\_delivered, BBR.target\_cwnd)$ ;
151            **else if** $cwnd < BBR.target\_cwnd$ **or** $BBR.delivered < InitialCwnd$ **then**
152                $cwnd = cwnd + packets\_delivered$ ;
153            $cwnd = max(cwnd, BBRMinPipeCwnd)$ ;
154        *BBR_Modulate_Cwnd_For_Probe_RTT()* ;

155    BBR_Update_Target_Cwnd():
156        $BBR.target\_cwnd = BBR\_Inflight(BBR.cwnd\_gain)$ ;

157    BBR_Inflight(gain):
158        **if** $BBR.RTprop == Inf$ **then**
159            **return** *InitialCwnd* ;                                          /* no valid RTT samples yet */
160        $quanta = 3 * BBR.send\_quantum$ ;
161        $estimated\_bdp = BBR.BtlBw * BBR.RTprop$ ;
162        **return** $gain * estimated\_bdp + quanta$ ;

163    BBR_Modulate_Cwnd_For_Recovery():
164        **if** $packets\_lost > 0$ **then**
165            $cwnd = max(cwnd - packets\_lost, 1)$ ;
166        **if** *BBR.packet_conservation* **then**
167            $cwnd = max(cwnd, packets\_in\_flight + packets\_delivered)$ ;

168    BBR_Modulate_Cwnd_For_ProbeRTT():
169        **if** $BBR.state == ProbeRTT$ **then**
170            $cwnd = min(cwnd, BBRMinPipeCwnd)$ ;

**Packet-transmission function definitions:**

171    BBR_Handle_Restart_From_Idle():
172        **if** $packets\_in\_flight == 0$ **and** *C.app_limited* **then**
173            $BBR.idle\_start = true$ ;
174            **if** $BBR.state == ProbeBW$ **then**
175                $BBR\_Set\_Pacing\_Rate\_With\_Gain(1)$ ;

**Save and Restore cwnd function definitions:**

176    BBR_Save_Cwnd():
177        **if not** *In_Loss_Recovery()* **and** $BBR.state! = ProbeRTT$ **then**
178            **return** *cwnd* ;
179        **else**
180            **return** $max(BBR.prior\_cwnd, cwnd)$ ;

181    BBR_Restore_Cwnd():
182        $cwnd = max(cwnd, BBR.prior\_cwnd)$ ;

---

**Algorithm 13:** Coupled Multipath BBR (C-MPBBR) - [theoretical part in 2.3.7]

---

**Fulfilling Goal 1 (MPTCP incentive):**

1   $\beta = 40$,   $Total\_Del\_Rt = 0$ ;

2   $lowest\_bw\_among\_all\_SFs = 999999999$ ;

3   $highest\_bw\_among\_all\_SFs = 0$ ;

4   $total\_number\_of\_SFs = 0$ ;

5   **if** $cmpbbr \rightarrow mode = CMPBBR\_PROBE\_BW$ **and** $cmpbbr \rightarrow cycle\_index = 3$ **then**

6    **for all** $SF_i$ **do**

7     **Calculate:** $Del\_Rt\_of\_SF_i$ ;

8     $total\_number\_of\_SF_s + +$ ;

9     $Total\_Del\_Rt + = Del\_Rt\_of\_SF_i$ ;

10     **if** $lowest\_bw\_among\_all\_SFs > bw\_of\_SF_i$ **then**

11      $lowest\_bw\_among\_all\_SFs = bw\_of\_SF_i$ ;

12     **if** $highest\_bw\_among\_all\_SFs < bw\_of\_SF_i$ **then**

13      $highest\_bw\_among\_all\_SFs = bw\_of\_SF_i$ ;

14   $threshold\_bw\_for\_stopping\_lowest\_bw\_SF = highest\_bw\_among\_all\_SFs * (1 - \frac{\beta}{100})$ ;

15   **if** $threshold\_bw\_for\_stopping\_lowest\_bw\_SF > Total\_Del\_Rt$ **and**
   $cmpbbr \rightarrow last\_number\_of\_SFs\_in\_btlneck < 2$ **and** $total\_number\_of\_SFs > 1$ **and**
   $lowest\_bw\_among\_all\_SFs \neq highest\_bw\_among\_all\_SFs$ **then**

16    $cmpbbr \rightarrow stop\_lowest\_bw\_SF\_count + +$ ;

17   **else**

18    $cmpbbr \rightarrow stop\_lowest\_bw\_SF\_count = 0$ ;

19   **if** $cmpbbr \rightarrow stop\_lowest\_bw\_SF\_count \geq 5$ **and** $total\_number\_of\_SFs > 1$ **and**
   $bw\_of\_SF_{this} = lowest\_bw\_among\_all\_SFs$ **then**

20    $cmpbbr \rightarrow stop\_lowest\_bw\_SF\_count = 5$ ;

21    **Close:** $SF_{this}$        `/* Close` $SF_{this}$ `whose` $Total\_Del\_Rt < \beta * highest\_bw\_among\_all\_SFs$
                             `for five successive ProbeBw states */`

**Fulfilling Goal 2 (MPTCP fairness):**

22   $\alpha = 20$,   $number\_of\_SFs\_in\_btlneck = 0$ ;

23   **if** $cmpbbr \rightarrow mode = CMPBBR\_PROBE\_BW$ **and** $cmpbbr \rightarrow cycle\_index = 3$ **then**

24    $bw\_lower\_limit = bw\_of\_SF_{this} * (1 - \frac{\alpha}{100})$ ;

25    $bw\_upper\_limit = bw\_of\_SF_{this} * (1 + \frac{\alpha}{100})$ ;

26    **for all** $SF_i$ **do**

27     **if** $bw\_of\_SF_i \geq bw\_lower\_limit$ **and** $bw\_of\_SF_i \leq bw\_upper\_limit$ **then**

28      $number\_of\_SFs\_in\_btlneck + +$ ;

29    **if** $number\_of\_SFs\_in\_btlneck > 1$ **and** $cmpbbr \rightarrow last\_number\_of\_SFs\_in\_btlneck > 1$ **then**

30     $final\_number\_of\_SFs\_in\_btlneck = number\_of\_SFs\_in\_btlneck$ ;

31    **else if** $number\_of\_SFs\_in\_btlneck = 1$ **and** $cmpbbr \rightarrow last\_number\_of\_SFs\_in\_btlneck > 1$ **then**

32     $final\_number\_of\_SFs\_in\_btlneck = cmpbbr \rightarrow last\_number\_of\_SFs\_in\_btlneck$ ;

33    **else**

34     $final\_number\_of\_SFs\_in\_btlneck = 1$ ;

35    $cmpbbr \rightarrow last\_number\_of\_SFs\_in\_btlneck = number\_of\_SFs\_in\_btlneck$ ;

36   $bw\_of\_SF_{this} = bw\_of\_SF_{this} / final\_number\_of\_SFs\_in\_btlneck$ ;