

TRANSPORT CONTROL PROTOCOL IN OPTICAL BURST SWITCHED NETWORKS: ISSUES, SOLUTIONS, AND CHALLENGES

BASEM SHIHADA AND PIN-HAN HO, UNIVERSITY OF WATERLOO

ABSTRACT

Since its advent in 1981, TCP has been subject to a tremendous amount of research effort and enhancements for achieving better performance over various network environments and application scenarios. Due to the transmission characteristics of optical burst switched networks, such as random burst dropping, retro-blocking (i.e., bursts proceeding or delayed from their actual reservation time slot), burstification delay, and burst signaling delay, TCP could be significantly affected if no corresponding countermeasure and enhancement are developed. In this review article we provide a comprehensive survey on reported studies for TCP enhancements over OBS networks in order to mitigate the numerous side effects due to the bufferless characteristic of burst transmission. Furthermore, we closely analyze TCP behavior over OBS networks with various burst transmission characteristics while highlighting the open challenges that have not yet been extensively tackled or solved.

Currently, the complex Internet infrastructure interconnects millions of communication devices and computing equipments through wired and wireless connections. One of the key success factors of today's Internet infrastructure has greatly relied on the ability to maintain a reliable, self-regulated, and congestion-tolerant transport protocol that serves end-user applications. Transport Control Protocol (TCP) [1], which was designed for military communication by ARPANET [2], has been taken to serve as the most pervasive and well recognized standard for a majority of currently available Internet-based applications. While analyzing the statistical network data, it has been reported that TCP is the most predominant protocol in terms of traffic volume (in bytes), which may take up to 90 percent of the total Internet traffic [3]. High-performance computing (HPC) networks observed an average share of 83 percent, and for the NETI@Home data the authors found that TCP flows contribute the major traffic volume seen in our datasets.

INTRODUCTION TO TCP

TCP congestion control mechanisms can be classified into the following three categories:

- Loss-based (e.g., Reno [4] and Sack [5])
- Delay-based (e.g., TCP Vegas [6] and Fast TCP [7, 8])
- Explicit notification-based (e.g., XCP [9])

These are basically the approaches taken by a TCP sender to determine if the network is in a congestion state such that the transmission rate is adjusted accordingly, while the receiver can be totally reactive to the transmission protocol.

The standard dropping-based TCP protocol stack, TCP Reno, which follows an additive increase multiplicative decrease (AIMD) window-based congestion control mechanism for regulating data transmission and maintaining network bandwidth usage, simply takes a packet (or TCP segment) dropping event as an indication of network congestion. The AIMD framework is composed of four stages for each TCP sender: *slow start*, *congestion avoidance*, *fast transmission*, and *fast recovery* [4].

At the slow start phase, the TCP sender has just started or restarted transmitting due to a timeout (TO) event, and the congestion window (*cwnd*) is initialized to the size of one segment. The size of *cwnd* is exponentially increased at each successfully delivered and acknowledged segment. The exponential growth of the *cwnd* continues until either it exceeds the receiver's advertised *cwnd* or a packet loss is

reported, which leads the TCP into the congestion avoidance stage. Note that the network congestion state is reported when either the transmission TO threshold is reached or the receiver collects triple-duplicated acknowledgments (TD). In the former case where a TO event occurs, which indicates heavy congestion, the *cwnd* is set to one segment, and the TCP returns to the slow start stage followed by the congestion avoidance stage. In the latter case, which indicates light congestion, TCP enters the fast retransmission stage by taking half of the sender's *cwnd* as the slow start threshold *ssthresh* and setting the *cwnd* to *ssthresh* plus three segments. Then the sender retransmits the lost segment and increments its *cwnd* by the segment. After acknowledging the second data segment, the *cwnd* is set to the *ssthresh*.

Different from dropping-based TCP, delay-based TCP implementations, such as TCP Vegas [6] and Fast TCP [7, 8], estimate the bandwidth and congestion status in the network by measuring the delay of each packet transmission in terms of round-trip time (RTT). The performance of Fast TCP has been evaluated in [7, 8]; it can be considered as an enhancement of TCP Vegas with better throughput. Nonetheless, in the following context we focus on introducing the TCP Vegas congestion control scheme for better comprehensiveness.

TCP Vegas modifies TCP Reno in the congestion avoidance, slow start, and retransmission stages, and determines the congestion status in the network according to the comparison of estimated and measured throughputs. TCP Vegas first computes *BaseRTT* as the minimum measured RTT, which is mainly determined by the propagation delay and queuing delay. The *Expected* throughput can thus be derived as *Expected* = *cwnd*/*BaseRTT*, where *cwnd* is the current congestion window size. Second, the *Actual* throughput is calculated for every RTT using the most recent measured RTT: *Actual* = *cwnd*/*RTT*. The TCP Vegas sender then computes the difference between the two quantities (denoted *Diff*):

$$\begin{aligned} Diff &= (Expected - Actual)BasedRTT, \text{ where } Diff > 0 \\ &= \left(\frac{cwnd}{BaseRTT} - \frac{cwnd}{RTT} \right) BaseRTT \\ &= cwnd \left(1 - \frac{BaseRTT}{RTT} \right). \end{aligned}$$

Obviously, *Diff* is non-negative, and is used to adjust the next *cwnd*. Vegas defines two threshold values, denoted α and β , for controlling the *cwnd*, which is changed as follows:

$$cwnd = \begin{cases} cwnd + 1 & diff < \beta \\ cwnd & \alpha \leq diff \leq \beta \\ cwnd - 1 & diff > \beta \end{cases}$$

If $Diff < \alpha$, Vegas increases *cwnd* linearly in the next round. If $Diff < \beta$, Vegas decreases *cwnd* linearly in the next round. Otherwise, Vegas leaves the *cwnd* unchanged.

The TCP Vegas congestion avoidance mechanism aims to maintain the expected number of backlog packets on the fly in the network between α and β . If the *Actual* throughput is much smaller than the *Expected* throughput, it is likely that the network is congested, and the TCP sender should reduce the flow transmission rate. On the other hand, if the *Actual* throughput is close to the *Expected* throughput, the connection may not be utilizing the available bandwidth, so the flow rate should be increased.

In the slow start stage, TCP Vegas increases *cwnd* exponentially only every other RTT, and exits slow start if *cwnd*

reaches the slow start threshold. In between two consecutive RTTs, the *cwnd* stays fixed in order to make a valid comparison between the *Expected* and *Actual* throughput.

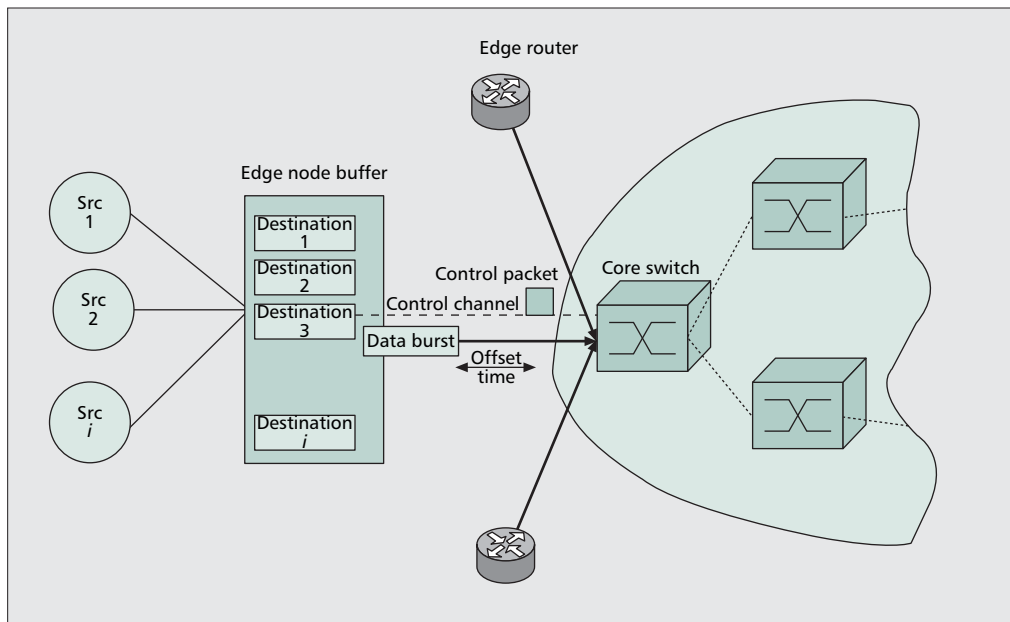
When TCP Vegas sender receives an acknowledgment (*ack*), it records the clock and calculates the estimated RTT using the current time and the timestamp recorded for the associated packet. Vegas then decides whether to retransmit the packet based on the following two conditions: First, when a duplicate *ack* is received, Vegas checks if the difference between the current time and the timestamp recorded for the associated packet is greater than the TO value. If true, the Vegas sender retransmits the packet without having to wait for the remaining incoming duplicate *acks*. Second, when an *ack* is received, if it is the first or second *ack* after a retransmission, Vegas again checks if the time interval since the packet was sent is larger than the TO value. If it is, the Vegas sender retransmits the packet. This will catch any other packet that may have been lost prior to retransmission without having to wait for the remaining duplicated *acks*; hence, the Vegas retransmission mechanism reduces the time to detect a lost packet from the third duplicate *ack* to the first or second one. After the packet retransmission is triggered by a duplicate *ack*, the *cwnd* is reduced by 25 percent (instead of 50 percent in Reno) only if the time since the last window size reduction is more than the current RTT.

The third TCP congestion control mechanism belongs to explicit notification-based TCP. In this category TCP with Explicit Loss Notification (ELN) [10] or Explicit Congestion Notification (ECN) [11] are the two representative approaches, which were developed to identify suspicious packet losses that may cause unnecessary retransmissions. Furthermore, TCP ELN/ECN can distinguish packet losses among congestion, contention, link failure, or other reasons. Packet retransmission starts as soon as the ECN or ELN is received, which is the round after the one that encounters a packet loss.

EVOLUTION OF OPTICAL SWITCHING TECHNOLOGIES

Optical backbones based on wavelength-division multiplexing (WDM) are the most prevalent transportation carrier in modern communication networks. Several switching technologies have been proposed to take advantage of the high-transmission capacity of fiber optics. Early approaches followed optical circuit switching (OCS), where point-to-point lightpaths are established for a relatively long period of time. OCS follows the store and forward approach where each optical switch performs optical-to-electrical-to-optical (O/E/O) conversions. However, the rapid increase in user demands and traffic engineering requirements have imposed several limitations on the adoption of OCS. First, since the number of available wavelengths is limited to a few tens or hundreds, the task of routing and wavelength assignment (RWA), which optimizes the assignment of wavelengths to all users, is NP-hard. Second, due to its quasi-static nature, OCS can barely support dynamic traffic variation and frequent connection requests.

Due to the ubiquity of the Internet Protocol (IP), much research has addressed the integration of IP with WDM networks. Optical packet switching (OPS) [12–14] has been proposed to support this integration, where the optical core can be taken as an extension of the IP layer. With OPS, each optical packet consists of a header and data payload, and is launched into the optical network. As a packet arrives at an optical switch, the header is converted and read in the electronic domain to configure the switch fabric, while the data is buffered optically in a fiber delay line (FDL) [15] until the switch configuration is completed. OPS follows an all-optical (AO) buffer-oriented transmission approach. OPS can suc-



■ **Figure 1.** An illustration of OBS networks.

cessfully solve the inadequacy and inefficiency of the OCS technology in terms of bandwidth provisioning, dynamics, and capacity utilization. However, the technical barriers before such a highly synchronized system can be commercialized are the high cost of all-optical buffer facilities and the dimensioning of these delay lines.

Optical burst switched (OBS) networks have attracted the attention of researchers due to their ability to achieve more dynamic and on-demand bandwidth allocation than OCS. They also offer improved network economy, and enable control and management integration [12, 13]. Compared to OPS, OBS is more practical to implement, and combines the best of OCS and OPS networks. With OBS, a data burst is formed at the edge node by assembling multiple incoming packets with the same destination and/or quality of service (QoS) requirements. To transfer the burst to the destination, a corresponding control packet that contains both the burst size and the burst arrival time is created at the network edge node, and sent prior to the launch of the burst [12, 13]. Since the bursts cut through each intermediate node in the network core while the control packets are subject to processing at each core node, a certain amount of offset time must be imposed between launching control packets and the corresponding bursts. The calculation of the offset time has to consider the upper bound on the number of hops and nodal processing delay. OBS follows the AO transmission approach. With its out-of-band signaling mechanism, OBS provides complete separation of the control and data domains, which can yield better network manageability and flexibility. Since the launched bursts cut through the network core without any buffering, the bursts can be subject to a minimum amount of delay.

OPTICAL BURST SWITCHED (OBS) NETWORKS

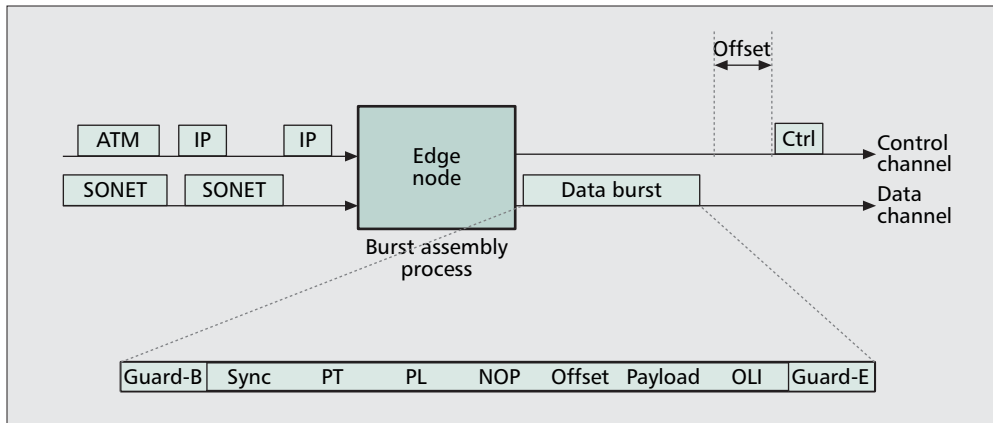
OBS networks are basically bufferless yet best effort in nature [12, 13, 16–18]. Basically, with IP over OBS networks, the IP applications and services are overlaid on the OBS backbone, and the two layers are interconnected through OBS edge nodes and the corresponding IP routers. The term *OBS domain* generally depicts the optical network with WDM technology formed by a group of OBS edge and core nodes, which run a suite of dedicated routing and signaling protocols. The

OBS domain routing and signaling protocols could be independent of the upper IP-compatible protocols, where an overlay model is followed. Alternatively, the IP and OBS domains may run a joint protocol by making the whole network into an augmented service model, where some extent of routing and traffic engineering information exchange is performed. Our discussions in this article focus on the latter case in order to achieve a cross-layer design.

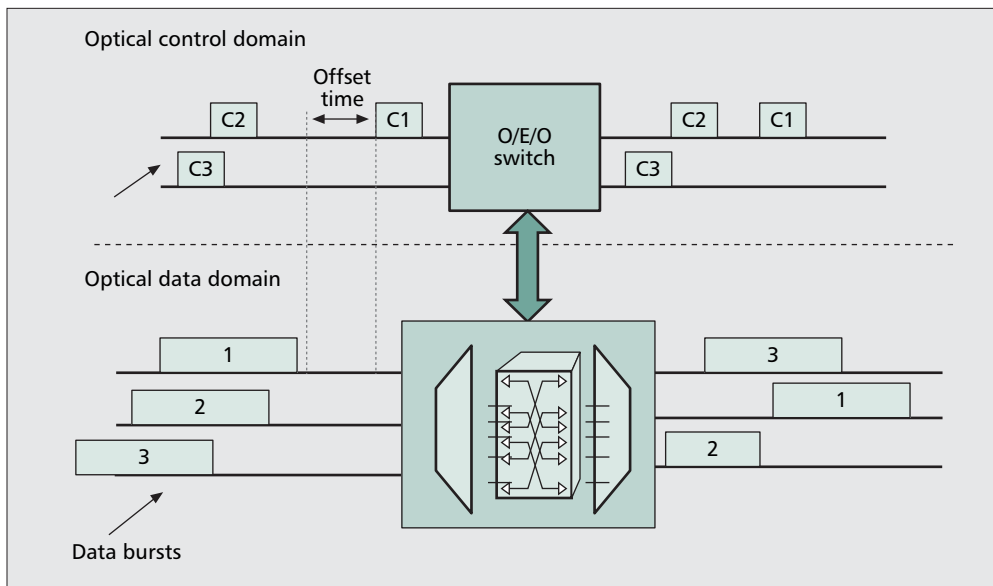
Within the OBS domain, a data burst is formed at an OBS ingress edge node by assembling multiple incoming IP packets with a common OBS egress node and similar QoS requirements [12, 16], as shown in Fig. 1. Since each IP-access router is attached to an OBS edge node, the OBS domain which provides high-speed transportation, can serve as the link-layer transmission media for the IP-access routers.

The bufferless and AO nature make the OBS networks distinguished from the traditional packet-switched networks which significantly rely on the electronic processing capability [16]. The implementation of OBS networks requires more precise signaling and a much higher switching speed than required by OCS networks [17, 18]. The edge nodes are capable of performing electronic-optical (E/O) and optical-electrical (O/E) conversion to collect/sort/assemble incoming IP packets from the higher layers into optical data bursts and disassemble the optical data bursts back to the IP packets at the destination edge nodes. The edge nodes are also responsible for signaling and routing in the OBS domain. The core switches, on the other hand, are subject to configuration requests via control packets sent by the OBS edge nodes, where the data bursts are forwarded and cut through the OBS domain in an AO manner [12, 13, 16–18].

A burst contains several other fields when assembled at the edge node. Figure 2 illustrates the burst assembly process as well as the burst format. A burst consists of guard preamble (Guard-B) and guard postamble (Guard-E) fields that help to overcome the uncertainty of data burst arrival time and data burst duration due to clock drifts between different core nodes. The payload basically consists of the assembled incoming data packets. Payload type (PT), payload length (PL), and number of packets (NOP) represent the type of data, data length, and number of packets assembled in the burst, respectively. The offset of padding indicates the first byte of padding. The optical layer information (OLI) includes



■ Figure 2. OBS edge node architecture.



■ Figure 3. OBS core switch architecture.

some information for performance monitoring and forward error correction obtained from the communication channels [12].

Typically, an OBS core node (switch) consists of two layers. The upper layer is responsible for processing control packets and configuring the switching fabric. In this layer control packets are processed, switching resources are reserved, and switching resources are freed after the burst cut through the switch. The switch matrix control unit and the port forwarding table, which is a lookup table and link scheduling [19–23] module are maintained. The lower layer is responsible for all-optical burst transport functionality. The lower layer consists of optical ports, wavelengths, and optical-to-optical connections. Figure 3 illustrates the generic OBS core switch architecture.

Fixed-path routing is a general approach taken in the OBS domain, where a physical route is predefined for each pair of OBS edge nodes. In order to deliver a burst from the source OBS edge to a certain remote one, a corresponding control packet that keeps the length of the burst and the arrival time of the burst at each intermediate node of the given route is created at the edge node and is sent prior to the launching of the data burst [17]. Since the burst cuts through the intermediate core nodes without any processing delay, the control packet is intrinsically slower than the burst. In order to pre-

vent the burst from catching up with its control packet, it is important to define offset time between the time instants of launching the burst and the control packet [17, 18], as shown in Figs. 2 and 3. This is also referred to as *delayed reservation*, which serves as the basis for OBS operation. The offset time is added to the average delay of the burst delivery with its minimum value determined by the number of hops along the route that the control packet traverses through. Then the data burst cuts-through the network following their control packet in a one-way (i.e., Tell-and-Go) signaling fashion.

Although OBS can achieve much better flexibility and efficiency than OCS-based wavelength-routed networks, OBS suffers from burst contention, which leads to burst drop due to the one-way resource reservation signaling and the lack of optical buffering. Burst contention occurs when more than one burst simultaneously attempt to take a common output port or wavelength channel [12, 16], which results in dropping $n - 1$ bursts when there are n contended bursts. Burst retransmission [24, 25], burst deflection routing [26, 27], fiber optic delay lines (FDL) [15], burst segmentation [28], and wavelength conversions [29] are some reported approaches to increasing the reliability of the burst transmission in the OBS domain. Each scheme has some impact on TCP performance, which has motivated numerous related studies to be seen in the past decade.

OVERVIEW OF THE ARTICLE

The goal of this article is to give an overview of the state of the art of TCP developments over OBS networks, including the recent proposed modifications and research challenges in such a deployment scenario. In addition, we are committed to conducting a comprehensive survey and tutorial on the reported solutions in TCP over OBS networks by classifying them into the following three categories:

- Link layer solutions (i.e., solutions implemented at the OBS domain)
- Modifications or enhancements based on conventional TCP (TCP originally designed for IP packet-switched networks)
- Newly designed TCP congestion control mechanisms

In the first category the solutions, which aim to assist TCP in understanding various burst transmission behaviors, require an implementation of additional supporting functions at either the OBS edge node or the OBS core node. Generally there are two problems incurred in this solution category. First, it breaks the TCP end-to-end semantics. Second, it introduces additional control overhead that impairs the applicability of the scheme. The solutions in the second category can overcome the above mentioned problems by preserving the end-to-end semantics of TCP congestion control such that no extra signaling overhead is required by any intermediate network device except the TCP senders and receivers. However, the price is paid in the effectiveness of the scheme, which may make it hard for TCP senders to follow the fluctuation of network congestion state in the OBS domain. Finally, the third category of solutions includes a totally new transport protocol that copes with the underlying burst transport behaviors.

The rest of the article is organized as follows. The issues of running TCP over OBS networks are identified, where the problems due to the bufferless transmission characteristics in OBS networks are formulated. A comprehensive survey is given of the recently proposed solutions for TCP over OBS networks in all three categories. Further discussions on the algorithms and some open research problems are given. Finally, we summarize and conclude the article.

ISSUES OF TCP OVER OBS NETWORKS

In this section the taxonomy of OBS networks, barebone OBS and burst contention resolution (BCR) OBS, is introduced. The transmission characteristics of each type are further identified. Afterward, we discuss the fundamental problems of running TCP over both types of OBS networks.

TAXONOMY OF OBS NETWORKS

Considering the OBS resource reservation signaling mechanisms, implicit [17] and explicit [18] resource reservations have been reported. In terms of burst assembly mechanisms, time-based, burst-length-based, mixed (both time and burst length), and optimized burst-size-based OBS networks have been reported [12, 13, 16]. Since burst contention resolution mechanisms play a vital role in the deployment of TCP over OBS, this article classifies OBS networks according to whether there are any burst contention resolution mechanism in the OBS domain, where two categories can be derived: barebone OBS and BCR OBS.

Barebone OBS — In barebone OBS a control packet is sent to reserve the intermediate switching fabric in a one-way signaling scheme. After a predefined offset time, the burst cuts through a certain path until it reaches the destination. In the

literature two major resource-reservation protocols in OBS domain — Just-In-Time (JIT) and Just-Enough-Time (JET) — are widely investigated. The JIT resource reservation protocol introduced in [18] follows the explicit reservation approach. Upon receiving the control packet, the switching fabric is reserved. The control message continues traveling on a hop-by-hop basis followed by the data burst until the destination is reached. Once the burst leaves the edge node, a teardown packet is sent along the same route to release the reserved resources.

The JET resource reservation protocol introduced in [17] follows the delayed reservation with an implicit resource release approach. The intermediate core nodes traversed by the burst can automatically release the reserved switching fabric shortly after the burst leaves the switch. Since both schemes never guarantee channel availability, if two or more data bursts attempt to access a common output port simultaneously, one of them must be dropped. Dropped bursts will be lost since barebone OBS does not employ any buffering or BCR mechanism. It is worth noting that the burst contention losses occur even while the traffic load is low.

OBS with Burst Contention Resolution — In contrast to barebone OBS networks, BCR OBS defines efforts in burst contention resolution such as retransmission of contended bursts at the edges [24, 25], burst deflection routing [26, 27], optical buffering [15], burst segmentation [28], and wavelength conversion [29]. With burst retransmission, deflection, or buffering, bursts subject to contention in the OBS domain can be retransmitted at the edge node, deflected to an alternate route at the core node, or buffered using FDL at the core node, respectively. Hence, burst random contention, in which a data burst is dropped when the channel utilization is low, may contribute a significant portion of the total burst loss probability. With burst segmentation, the contended bursts are segmented at the core node such that overlapping segments are discarded while the rest of the burst can be transmitted successfully. With wavelength conversion, optical data flows on a certain wavelength can be converted to another wavelength, which can effectively reduce the resource segmentation and result in a better chance of successfully forming a data path.

The above BCR schemes have been proven to effectively increase transmission reliability and throughput in the OBS network domain at the expense of introducing extra overhead, delay, and computation complexity, which are discussed in detail later.

CHARACTERISTICS OF OBS NETWORKS

The bufferless nature of OBS results in two major transmission characteristics that distinguish OBS networks from any other type of network:

- Random burst contention losses at low traffic loads
- Assembling (burstification) and signaling delay

Assembling (Burstification) and Signaling Delay — Several burst assembly algorithms have been proposed for OBS networks and can be divided into the following four categories [12, 13, 16]. The first category is time-based, in which a timer is set after the beginning of each new assembly cycle. With this assembly algorithm, all IP packets arriving within a specific period of time are considered to be assembled into a burst assuming they have a common destination edge node. The second set of algorithms is burst-length-based, where a burst length threshold is set. With these schemes, the burst length threshold serves as the minimum burst length before the burst

leaves the network. The third category combines time- and burst-length-based, where bursts are assembled and sent when either the burst length exceeds the desirable threshold or the timer expires. The last category includes adaptive assembly algorithms, where a dynamic adaptive threshold on the burst length is set in order to optimize the overall performance in the OBS domain for QoS sensitive traffic [14]. Both barebone and BCR OBS experience burstification and burst signaling delays.

Random Burst Contention Losses — In IP-based networks with electronic processing, contended packets can be temporarily stored in the memory of intermediate IP routers instead of being dropped immediately. Except for extreme cases such as buffer overflow caused by persistent network congestion, packets can be delivered in a hop-by-hop manner until they reach the destination. Nonetheless, this is not the case in OBS networks. As a compromise between OCS and OPS networks, OBS facilitates statistical multiplexing to some extent by using JIT or JET resource reservation schemes while buffering the bursts at the OBS edges only. Due to the adoption of a one-way signaling mechanism and the lack of global scheduling, burst contention and resultant dropping could happen even when the network is lightly loaded [16]. This is referred to as random burst contention losses, which may impose a significant vicious impact on TCP protocol stacks, especially those that take packet loss events as the only indication of network congestion. In contrast with random burst contention losses, the network could be subject to persistent network congestion, where a high utilization state of the OBS network lasts for a much longer time.

ISSUES OF TCP OVER OBS

It has been shown that the burstification process at OBS edge nodes influences the performance of TCP Reno [30]. TCP flows are classified into three categories: fast, medium, and slow. In the fast TCP flow category, TCP flow is subject to high IP access bandwidth such that the entire transmitted segments of the *cwnd* are assembled in a single burst. On the other hand, the medium and slow flows have a proportion of TCP segments in the *cwnd* assembled in a single burst due to limited IP access bandwidth. For medium to slow flows, a performance penalty is introduced on the TCP sending rate due to assembly delay, while the aggregation of multiple TCP segments due to burst assembly has given TCP a throughput performance gain [30, 31]. When a burst is lost due to either persistent network congestion or random burst contention loss, fast TCP flows return to slow start since the TCP segments of the entire *cwnd* are lost, which indicates severe network congestion [32]. However, the medium and slow flows enter the fast retransmission stage since some of the TCP segments of the current *cwnd* are lost [32].

Multiple TCP Segments Loss — In TCP over OBS networks, TCP performance is significantly affected by random burst contention loss [30–32], which causes unnecessary *cwnd* cuts. Since an OBS edge may correspond to millions of TCP senders [3], hundreds or thousands of TCP segments could be assembled and transported together in a single burst. At the occurrence of a burst loss, TCP with fast flows will be significantly affected by falling into false congestion detection, also called false TO (FTO) [32]. However, in the case of medium and slow flows, the probability of having multiple segments of a single TCP session assembled in a single data burst is low, a burst loss will less influence a single TCP flow. However, it results in affecting a large number of TCP flows that have a

segment assembled in the lost burst. Depending on the flow speed [30], TCP usually suffers from FTO [32] or performs unnecessary *cwnd* cuts [30, 32].

Packet Reordering — Out-of-order burst delivery may occur due to the burst assembly mechanism, where the TCP segments of a certain TCP session are assembled in two or more bursts that are sent in a different order from that of the TCP segments. Also, bursts may be sent in correct order but received out of order due to burst drop, retransmission, or delays from deflection. The frequency of out-of-order delivery depends on the time and burst length thresholds, the traffic shape, or the QoS criteria defined in the burst assembly algorithm. Out-of-order delivery may cause a TCP sender to improperly retransmit a TCP segment since the TCP receiver will simply issue TD *acks* by assuming that there is a sequence of missing segments. Depending on the number of packets assembled in a single burst, TD *acks* can take place in a single round. The frequency of the false TD *acks* problem is expected to dramatically affect the sender in the presence of a large number of TCP segments assembled in a single burst. The higher the rate at which TD *acks* occur, the smaller the TCP sender *cwnd* remains, which significantly throttles the throughput. This problem is also referred to as the global synchronization problem.

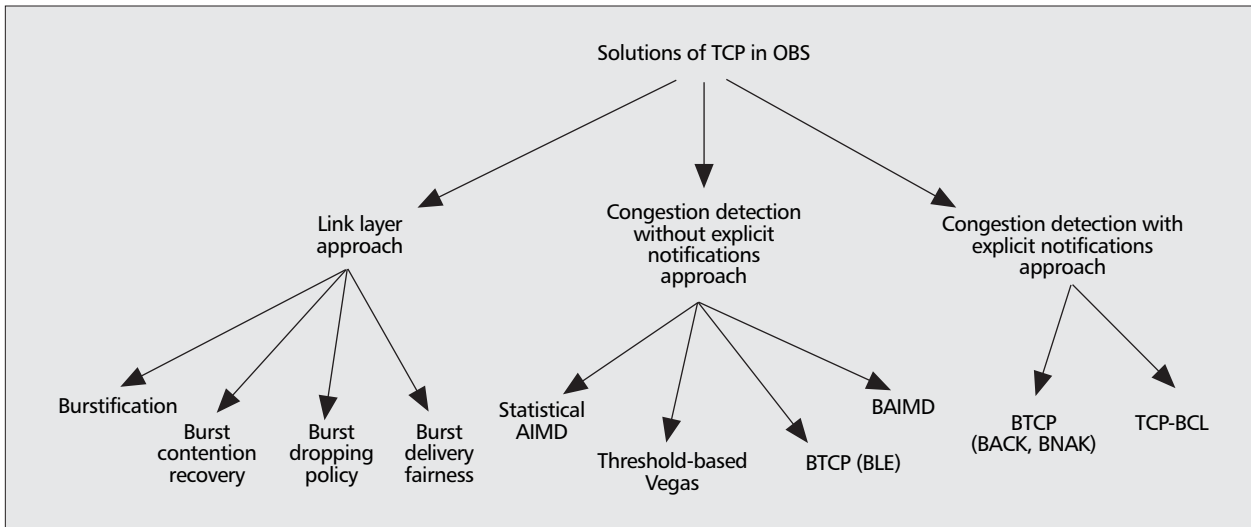
Slow Convergence — For high-bandwidth TCP flows that last for a long period of time and span over long distances, the TCP throughput may suffer from a very long convergence time due to the slow additive-increase congestion avoidance mechanism. The following gives an example. Let us have a 10 Gb/s TCP flow with 1.5 kbyte TCP packet size. Let the RTT be 100 ms. Thus, it takes about an hour to converge to the rate of 10 Gb/s assuming that the packet loss rate is less than 10^{-8} . The slow convergence problem is expected to be worse in OBS networks due to a longer RTT of each data burst caused by the burstification process and BCL delay [33].

With the additive increase of *cwnd* for each successful round of burst delivery, TCP over OBS networks may take a significant amount of time to converge to the available high bandwidth. On the other hand, the multiplicative decrease for each segment loss at the TCP senders is a critical response to network congestion, even if we exclude the possibility of any burst drop due to random burst contention. Note that it is necessary to maintain a large *cwnd* to achieve high throughput, which can only happen when an extremely low packet loss rate is seen. This is not likely to happen in the presence of random burst losses in OBS networks.

TAXONOMY OF SOLUTIONS FOR TCP OVER OBS

In bufferless OBS networks, TCP senders should not blindly consider the loss of data segments as an indication of network congestion; instead, TCP should attempt to collect further information and network states to distinguish between a packet loss event due to random burst contention and one due to persistent network congestion. TCP implementations of this category may require extra signaling efforts between TCP senders and the OBS domain [32], and/or within the OBS domain [24, 25, 34]. In this article three categories of previously reported solutions for improving TCP over OBS are briefly described as follows.

Link Layer Solutions — With link layer solutions, the OBS domain is simply treated as the link layer underneath the IP edge routers. Thus, OBS burst transmission is transparent to the upper TCP senders, where the TCP throughput is



■ Figure 4. The taxonomy of solutions for TCP in OBS networks.

improved by adopting mechanisms such as an adaptive burstification algorithm [14], a forward error correction (FEC) mechanism, BCR schemes, and automatic repeat request (ARQ) [24, 25]. The fundamental advantage of employing a link layer solution is hiding non-congestion burst loss from TCP senders such that the layered structure of network protocols is followed by keeping local burst transmission reliability. To achieve this link layer transmission reliability, burst retransmission, deflection, and adoption of FDLs in the OBS domain are common approaches. These approaches have contributed to reducing *cwnd* fluctuation and consequently improving TCP throughput at the expense of introducing additional delay and control overhead. In the rest of the article we refer to the schemes functioning only within the OBS domain as *link layer solutions* to reflect the fact that the OBS domain can simply be taken as the link layer transmission seen by the IP routers.

Congestion Detection without Explicit Notifications — With the solutions for congestion detection without explicit notifications, TCP senders try to evaluate the states in the OBS domain by way of some equations, measured RTT, and/or statistic tools. In this case TCP can maintain and analyze a number of previous RTTs at TCP senders in order to identify if a packet loss event is due to persistent congestion or random burst contention. Under this category, we identify the following sender-side congestion control schemes: Statistical AIMD, Threshold-Based Vegas, Burst TCP with burst length estimation (BLE), and OBS with general AIMD (BAIMD). There are other TCP solutions such as TCP Sack, duplicated Sack, and Fack that can work over OBS networks without explicit signaling [31, 32]. These schemes are excluded from our review since they were not originally proposed for OBS environments. However, the throughput performance of these TCP implementations over OBS networks was examined through simulation in [35].

In TCP fast flows a burst that is successfully retransmitted, deflected, or buffered is subject to additional delay due to the contention resolution mechanisms. The additional delay is inevitably added on top of the RTT of all the TCP segments assembled within the burst. In this case the delay-based TCP sender (e.g., TCP Vegas) detects a sudden increase in RTT of those segments, and interprets the additional delay as an indication of network congestion. This is another form of false congestion detection by delay-based TCP, which has been identified to fatally impair the TCP throughput due to unnecessary TO retransmissions followed by slow start at the TCP

senders. Threshold-based TCP Vegas was proposed to cope with the above mentioned problem [36, 37]. This scheme is considered a sender side TCP implementation that does not require any explicit notifications.

Congestion Detection with Explicit Notifications — The solutions of congestion detection with explicit notification are basically cross-layer designed such that the states in the OBS domain can be leaked to the upper TCP senders. In this case TCP senders can additionally be informed of the channel conditions and the actual cause of the packet loss, such as network congestion, packet corruption, or any possible hardware component failure. TCP senders will retransmit the lost segments without affecting their *cwnd* for packet losses due to any reason other than network congestion. This category folds the TCP implementations with explicit congestion or loss notifications such as OBS TCP (BTCP) (Back/BNack) and TCP with explicit burst contention loss (TCP-BCL).

As a summary of this section, we have briefly introduced and classified the previously reported solutions for TCP in OBS networks. As shown in Fig. 4, the following three main categories were highlighted: link layer solutions, TCP congestion detection without explicit notifications, and TCP congestion detection with explicit notifications.

OVERVIEW OF EXISTING SOLUTIONS

This section explicitly surveys all three categories in the aspects of dropping-based (e.g., Reno), delay-based (e.g., Vegas), and explicit-notification-based (e.g., TCP ELN/ECN) TCP implementations.

LINK LAYER SOLUTIONS

Solutions Based on Burstification Processes — It has been reported that the delay caused by the burstification process at the OBS ingress nodes will enlarge the TCP RTTs and may lead to TO if the TO threshold is not well manipulated, which has the potential to impair TCP throughput due to non-congestion-caused *cwnd* cuts [14, 30, 31]. In [30] the authors examined the performance of TCP Reno in the presence of various burst assembly algorithms at OBS ingress nodes. The authors have shown that TCP with medium to slow speed suffers from significant performance degradation due to the fact that the assembly delay becomes very large with respect to the arrival of each TCP segment. On the other hand, the authors

have shown that fast TCP flows are less influenced by the burstification delay.

The authors in [14] proposed an adaptive assembly period (AAP) algorithm and investigated the impact of some configuration parameters, such as burst assembly time, burst size, and threshold on adaptive assembly queue length, on the performance of both TCP and UDP traffic over OBS networks. The authors showed that the adaptive assembly algorithm supports achieving the best TCP performance among all the investigated assembly schemes.

The study in [38] has taken TCP Reno as a dominant TCP, where various burst assembly delays and burst sizes have been examined and simulated. The study concluded that TCP Reno has failed to deal with burst losses where each burst contains a large number of TCP segments assembled from a single TCP source. The study in [38] contributed with a TCP over OBS simulation package for NS-2, which was later extended by other research institutions [36–40].

Burst Contention Loss Recovery — As one of the link layer solutions, BCR schemes can be employed in the OBS domain in order to reduce random burst loss, thereby improving the transmission reliability of OBS networks. In many cases the reported BCR schemes have successfully hidden burst loss events from upper TCP senders at the expense of introducing extra transmission delay for bursts that experience contention. Among the three approaches, FDL facilitates the BCR by providing very limited buffering time for each contending burst. In general, 1 ms of buffering time requires a fiber with a length of 200 km. Thus, the resultant delay could be limited as well. On the other hand, burst retransmission and deflection may introduce a significant amount of delay. Therefore, in this survey study we focus on both burst retransmission and deflection.

With burst retransmission, the OBS edge node stores a copy of the launched bursts for possible retransmission. As the control packet traverses through the core nodes, an intermediate node that fails to reserve the resource will send an explicit notification to the OBS edge in order to report the reservation failure. Upon receiving the notification, the OBS edge retransmits a duplicate of the requested contending burst preceded by the corresponding control packet. The retransmitted burst will certainly experience an extra retransmission delay, which is the time elapsed between the initial control packet transmission and the last notification receipt for the corresponding burst.

If the network is lightly loaded, the retransmission scheme has a good chance of successfully delivering the contending bursts without involving TCP retransmission mechanisms. The studies in [24, 25] have shown that the retransmission scheme can significantly reduce the burst loss probability from that using a barebone OBS network, especially at low traffic load. The authors have also shown that retransmitting lost bursts from ingress nodes can avoid TCP false congestion detection. However, when the network is heavily loaded, the retransmitted bursts may still get blocked and finally lead to TO at the TCP senders. In this case the retransmission persistence (or the maximum number of retransmissions for a single burst) and time threshold at the OBS edge nodes before stopping the retransmission are important issues subject to further research efforts.

With burst deflection, a data burst is routed through its primary path if there is no burst contention. In the event of contention at the core node, the burst will be dynamically rerouted and directed to an alternative path segment starting at the core node where the burst encounters a contention. Since the primary path is usually the shortest path, the data

bursts following the alternative path segment result in a longer propagation delay [26, 27]. The study in [25] investigated burst retransmission along with deflection routing, which showed that the deflection scheme can significantly reduce burst loss probability at low traffic loads. The same as in burst retransmission, the extra delay induced by the deflection will cause some vicious effects on delay-based TCP. When the network is heavily congested, deflection may worsen the situation by consuming more resources, which leads to a fatal impairment of TCP throughput.

TCP with Burst Acknowledgment — In [34] a TCP throughput analytical model is introduced by considering the burst acknowledgment mechanism implemented in the OBS domain. The authors proposed an error recovery mechanism for electronic buffering at the edge nodes in order to improve the TCP throughput at the expense of taking extra memory space at the edge node for buffering a copy of all bursts within a certain time window. Once the amount of additional memory at the edge nodes becomes very large, the proposed scheme could be subject to a problem in practical implementation. Furthermore, introducing burst delivery acknowledgment violates the fundamental OBS burst transmission semantic. Let the edge node switching capacity be c , the number of IP packets be k , the average assembling granularity be M , and the burst dropping probability be B_{dp} . The required extra buffering space is $c \times (RTT_{burst} + B_{dp} \times (RTT_{burst} + k \times M))$, where RTT_{burst} includes the burst assembly time and burst offset time. The scheme cannot prevent TCP senders from receiving TO indications, where the $cwnd$ could be decreased in response to the burst retransmission delay.

TCP Decoupling — In [41] a modified TCP decoupling approach is introduced. This approach monitors the burst contention probability at the OBS network bottleneck link by taking advantage of the TCP self-clocking property, where the burst sending rate is controlled through the arrival time of TCP decoupling management packets. In this scheme a virtual circuit (VC) is set up for each source-to-destination edge node pair in the OBS domain. The VC is controlled through the TCP congestion control located at the OBS edges such that the sending rate never exceeds the link capacity. The OBS edge node uses TCP *ack* packets to control the timing of burst sending. Considering the simulation parameters provided in [41], the authors have demonstrated an improvement of TCP throughput by avoiding unnecessary burst losses, where the overall link utilization is increased from 50 to 62 percent, and the packet dropping probability is decreased from 50 to 30 percent. However, this approach requires maintaining a record of launched bursts, burst launch time, and the corresponding TCP segments for each source and destination pair. This approach complicates the OBS edge node architecture and functionality by manipulating the TCP packets (through TCP agents) at the OBS network.

Retransmission-Count-Based Dropping Policy — In [42] and [43], a dropping policy, called Retransmission-Count based dropping policy, is introduced, which aims to improve TCP throughput. The basic idea of the dropping policy works at the OBS edge node by taking the number of burst retransmission attempts into consideration, where the bursts that have been less retransmitted are dropped. It is known that the more frequently burst retransmission takes place in the OBS domain, the less time remains for the TCP timeout to be triggered. The authors proposed to add a retransmission count (RC) field in the burst control packet with an initial value of 1. In the event of burst contention, the core node compares

Schemes	Suspicious TOs	Packet reordering	Multiple packet losses in a round
Reno (Eifel)	✓		
New Reno (Eifel)	✓		
Sack (Eifel)	✓		✓
DSack		✓	✓
Fack			✓
DCR	✓	✓	✓
F-RTO		✓	✓
BTCP	✓	✓	✓
BAIMD	✓	✓	✓

■ Table 1. TCP implementation performance over OBS.

the control headers of the contending bursts and drops the bursts with lower RC values. The bursts with larger RC values have higher priority for successful retransmission since their assembled TCP packets have already experienced relatively longer delay. Once a burst is dropped, the corresponding NAK along with a copy of the RC value is sent back to the OBS edge node, which sets the RC field for the dropped burst to $(RC_{in_NAK}+1)$ and retransmits the burst. The number of retransmission attempts follows a predefined retransmission policy. This study aims to increase the transmission chances of TCP packets that have experienced the longest time in the OBS domain. However, the authors have not addressed the TCP fairness factor in the proposed dropping policy. Also, the preemption of reserved resources by bursts with larger RC values may lead to negative influences on the dropped flows, which may have already launched the corresponding data burst.

We observe that all the abovementioned link scheduling or signaling algorithms have demonstrated reduced burst dropping probability, thus increasing the overall network throughput. However, they have introduced extra switching architecture complexity and additional burst delay.

SOLUTIONS WITH EXPLICIT NOTIFICATIONS

Burst TCP — The study in [32] investigated TCP false timeout detection due to random burst contention loss under a wide range of traffic loads. Two solutions were proposed that use explicit notification. In the first approach [32], called BTCP with burst acknowledgment (Back), each TCP packet is acknowledged by a TCP agent located at OBS edge nodes. This approach can effectively prevent TCP from detecting false TO; however, the end-to-end TCP semantics is violated since *acks* reach TCP senders before the actual completion of packet delivery. The second proposed approach in [32], called BTCP with burst negative acknowledgment (Nack), maintains a TCP agent at each OBS core node. Whenever a burst is dropped, the TCP agent disassembles the burst and sends a burst Nack (BNack) to the corresponding TCP sender. The missing segments and the network congestion state are explicitly exchanged between the TCP senders and the OBS core nodes. In general, how to reduce the extra control overhead and implementation complexity is a challenge when attempting to deploy the abovementioned solutions.

TCP with Burst Contention Loss Notifications

— The article in [35] addresses the issues in the design of various TCP flavors in the OBS environment. TCP Reno, New Reno, Sack, Duplicate Sack [44], Fack, Forward Retransmission Timeout Recovery (F-RTO) [45], Eifel [46], delayed congestion response (DCR) [47], BTCP [32], and BAIMD [39] have been examined in terms of throughput performance under a wide range of burst drop probabilities. Table 1 summarizes the experiment results. The measured performance evaluation of both DCR and F-RTO are based on TCP Sack.

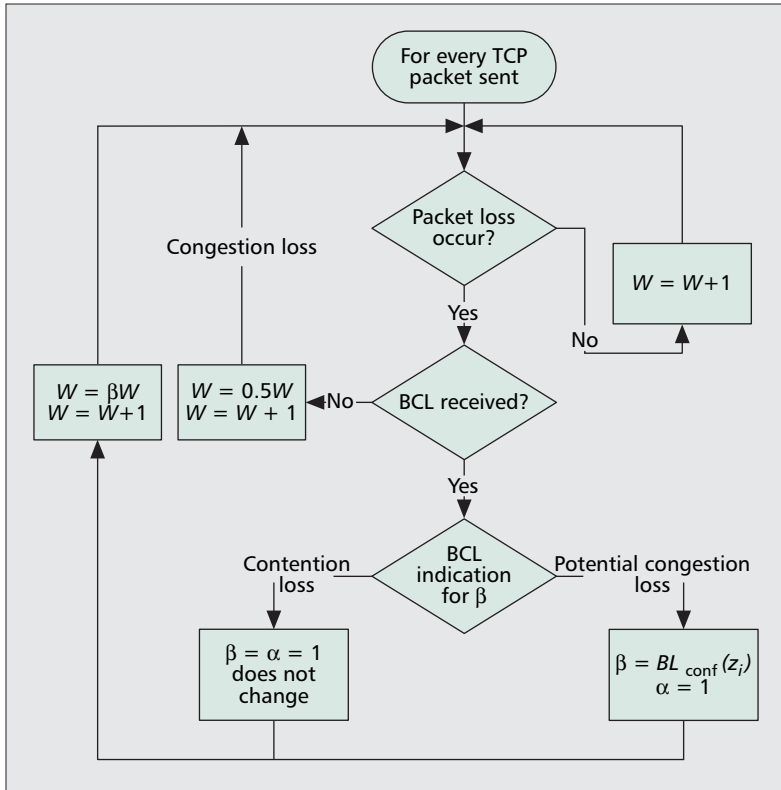
The simulation results have shown that the above TCP flavors fail to maintain a good throughput level in the presence of burst losses that contain multiple segments from a single TCP flow. Therefore, Explicit Burst Contention Loss Notification (BCL)-based TCP is introduced. The scheme aims to solve the false congestion detection problem in TCP over OBS networks and avoid the unnecessary *cwnd* reduction by adjusting TCP *cwnd* based on the utilization and burst dropping information carried in the explicit notification messages. This scheme is considered the first study that integrates the explicit notification platform with the GAIMD framework over OBS networks.

In this article the authors proposed a novel mechanism for detecting network congestion in bufferless OBS networks. In OBS networks congestion occurs at edge nodes when receiving packets at a much higher arrival rate (i.e., bursty impulse) than the edge node is capable of dealing with. This causes the edge node to drop bursts due to buffer overflow. The authors refer to the congestion at the network edge node as *edge congestion*. In addition to edge congestion, *path congestion* is defined, which leads to congestion in the network core nodes.

The authors proposed two possible approaches for detecting congestion at the network core (i.e., path). The first approach is to delegate the congestion detection process to the core nodes. The edge nodes receive explicit signals from the core switches indicating link congestion. A similar approach is proposed in BTCP with BNack [32]. It is notable that this approach may not be very practical since it increases the signaling and computation overhead at the core nodes. Therefore, the authors introduce a new mechanism for detecting congestion status along an OBS route (path) through the statistics gathered at the corresponding edge node. This approach does not introduce any additional signaling effort at the core nodes. Path congestion is measured based on how congested the route in the OBS domain is, which will be taken as an important index for the upper layer TCP senders using the route to adjust their congestion windows.

The authors proposed two possible approaches for detecting congestion at the network core (i.e., path). The first approach is to delegate the congestion detection process to the core nodes. The edge nodes receive explicit signals from the core switches indicating link congestion. A similar approach is proposed in BTCP with BNack [32]. It is notable that this approach may not be very practical since it increases the signaling and computation overhead at the core nodes. Therefore, the authors introduce a new mechanism for detecting congestion status along an OBS route (path) through the statistics gathered at the corresponding edge node. This approach does not introduce any additional signaling effort at the core nodes. Path congestion is measured based on how congested the route in the OBS domain is, which will be taken as an important index for the upper layer TCP senders using the route to adjust their congestion windows.

In this scheme each OBS ingress edge node maintains long-term and short-term statistics for each route initiated at it. Whenever a burst loss occurs, the ingress node determines if the route is congested by correlating the long-term and short-term statistics. In specific, let the parameter M be the number of transmitted burst along the OBS route, which is used to obtain the long-term statistics. Note that M should be sufficiently large to fully represent the intrinsic characteristics of the network topology, routing policy, and traffic pattern, etc. The outcome of the M burst deliveries is kept as a vector with each entry 0 or 1, which represents a burst drop event or successful delivery, respectively. Let the parameter N be the number of transmitted bursts for the short-term burst drop rate, which is generally small.



■ Figure 5. Flowchart of TCP-BCL over OBS.

The main idea of the proposed scheme is to position the average short-term burst drop rate (denoted as avg_b_N) in the spectrum of the long-term burst drop rate. To achieve this, the outcomes of the M burst deliveries are divided into

$$\left\lfloor \frac{M}{N} \right\rfloor$$

segments, each containing the outcomes of consecutive N burst deliveries. Thus, it is possible to obtain a vector denoted θ of a size

$$1 \times \left\lfloor \frac{M}{N} \right\rfloor,$$

where each entry, denoted θ_i for $i = 1$ to

$$\left\lfloor \frac{M}{N} \right\rfloor,$$

stores the burst drop rate of the i th small segment of burst deliveries. With this vector, it is possible to obtain the average burst drop rate in the M burst deliveries (denoted avg_b_M) and the variance of the vector (denoted var_b_M). If avg_b_N is much larger than avg_b_M , it is highly possible that the route in the OBS core is subject to random burst contention, and the corresponding TCP senders should not take it seriously in the response of $cwnd$ adjustment. On the other hand, with comparable avg_b_N and avg_b_M , we can expect that the current burst drop is more likely to be an indication of congestion along the route in the OBS domain.

Specifically, to quantify the relationship between avg_b_N and avg_b_M , the authors assume that the burst drop rate of each entry in θ can be positioned in histogram distribution.

Thus, it will be easy to position avg_b_N in the spectrum formed by the vector θ . A policy is created to define the relationship between avg_b_N and the spectrum formed by θ in order to define the confidence with which a burst loss event is due to network congestion is defined. A BCL is sent to the corresponding TCP senders if the burst loss is judged to be due to contention at low traffic load. Thus, TCP-BCL has the TCP senders take every segment loss event as due to congestion when a BCL is not received. Figure 5 illustrates the proposed TCP-BCL congestion control scheme. The proposed path congestion detection distinguishes the TCP-BCL scheme from the BTCP in [32] since TCP-BCL only signals the dropped bursts at low link utilization, while BTCP (Back/BNack) reports the loss of every burst loss from core nodes. Hence, the number of notifications in TCP-BCL will be much less than that in BTCP. Furthermore, this design can significantly reduce the intradomain signaling overhead as the edge node is the one responsible for detecting network congestion. This comes at the expense of introducing extra computation and signaling performed at the edge nodes.

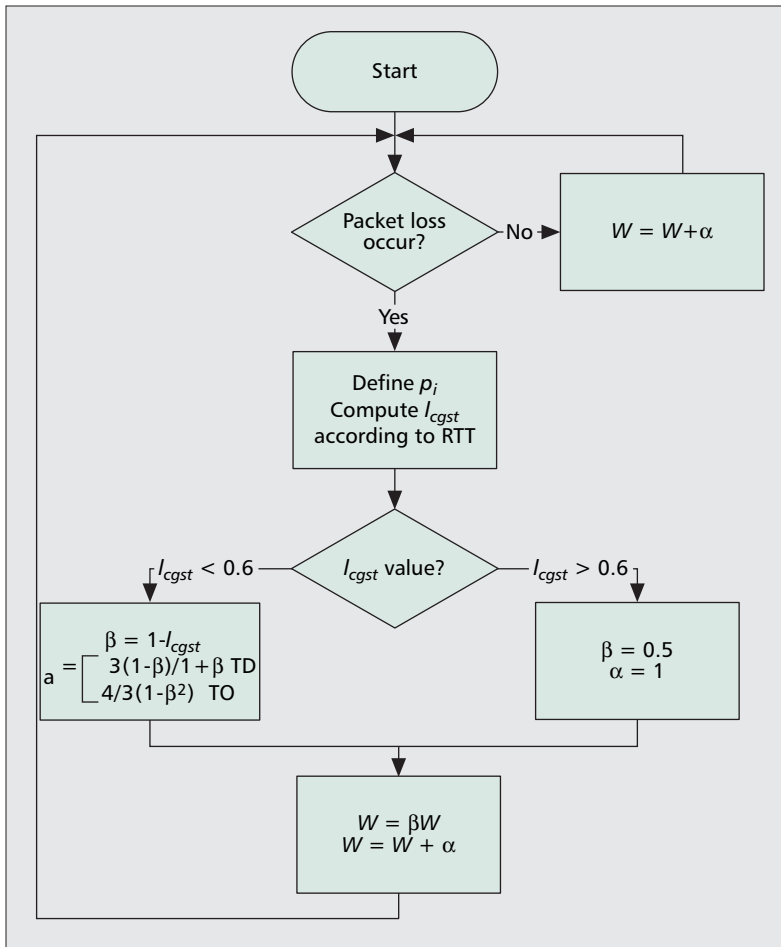
SOLUTIONS WITHOUT EXPLICIT NOTIFICATIONS

Burst TCP with Burst Length Estimation —

The study in [32] proposed solving false TCP TOs using a burst length estimation approach. This approach, called BTCP with Burst Length Estimation (BLE), is based on estimating the number of TCP packets assembled in the burst without the knowledge of the burst assembly algorithm deployed at the OBS edge node. In addition to the $cwnd$, a burst congestion window denoted $burst_wd$ is also maintained. When the TCP sender detects TO, it first compares its $cwnd$ with the $burst_wd$. If $cwnd \leq burst_wd$ and $burst_wd > 3$, the TCP sender considers this TO a false TO. It halves the $cwnd$ and performs fast retransmission for the missing segments. If $cwnd > burst_wd$ or $burst_wd \leq 3$, the TCP sender considers this a true TO and initiates a normal TCP retransmission procedure. This approach does not require any coordination or explicit information exchange between the TCP senders and the OBS network. However, the accuracy of the estimated $burst_wd$ remains an open research challenge. It is possible that a true TO can be taken as a false TO, which endangers the network's stability. Furthermore, this approach cannot distinguish the loss of multiple packets in a congested IP access network from the loss of multiple packets in a burst loss caused by random burst contention loss.

Burst AIMD (BAIMD) — The BAIMD [39] scheme is based on the framework of Generalized AIMD (GAIMD) [48, 49] for $cwnd$ adjustment. Two parameters are defined, α and β , which serve as the additive incremental and the reduction ratio for the $cwnd$ at each TCP sender. Unlike conventional TCP, where α and β are constantly set to be 1 and 0.5, BAIMD dynamically determines the two parameters in such a way that the $cwnd$ is increased by α segments for each acknowledged packet in a round and decreased by β ($0.5 < \beta < 1$) for any packet loss event. It is clear that BAIMD is more general than AIMD with much better design granularity and flexibility.

With BAIMD, the sender is not explicitly notified about the used burst assembly mechanism and the reasons for the



■ Figure 6. Flowchart of the BAIMD congestion control scheme.

burst losses. Each sender treats a packet loss event as a congestion loss. Obviously, this could lead to overestimation of network congestion by irrelevantly cutting the *cwnd* by half for every TCP segment drop. To compensate for this overestimation, BAIMD senders use β larger than 0.5. As such, the summarized effect of the burst drop event is estimated when both the values of α and β are dynamically determined in the BAIMD senders using burst-level states (i.e., the estimated traffic load in an OBS network). The scheme aims to achieve the best throughput for the competing flows. For example, if a burst is lost when the network load is low, the lost packets are considered due to random burst contention, and the multiplicative factor is set to be $0.5 < \beta < 1$. Otherwise, β is set to 0.5 when network load is heavy and burst dropping is due to congestion.

One of the most important advantages of BAIMD is being simple: no burst-level window is maintained at the TCP senders, and no explicit notification specific to each launched TCP segment is exchanged between the OBS edge and TCP senders. The scheme maintains clean separation between the control signaling at the TCP senders and OBS edge nodes. Most notably, BAIMD senders use the RTT of each launched segment along with the number of TOs as references for sensing the network load. For example, data bursts are subject to extra buffering delay at the OBS edge nodes in response to serious network congestion. Thus, RTT substantially increases. On the other hand, if the data burst is dropped due to random burst contention in barebone OBS, the RTT remains unchanged. Thus, it is considered an indication of low network load. Figure 6 illustrates the functionality of the BAIMD scheme.

BAIMD estimates the multiplicative factor β through estimating the traffic load at the OBS layer. BAIMD defines a certain threshold load value denoted l_{cgst} for computing the multiplicative factor. When a packet loss occurs at time t and all the connected TCPs as well as BAIMD are notified through either TD or TO, the maximum link capacity has been reached at that moment. In the congestion state (i.e., the used threshold is $l_{cgst} \geq 0.6$), BAIMD senders behave similar to conventional TCP senders with $\beta = 0.5$. Otherwise, β follows $\beta = 1 - l_{cgst}$ where, $0 < l_{cgst} < 0.6$. Once β is computed, BAIMD uses the GAIMD congestion control mechanism to obtain the sending rate α as follows. In the presence of TD loss, $\alpha = 3(1 - \beta)/(1 + \beta)$, while in the case of TO loss,

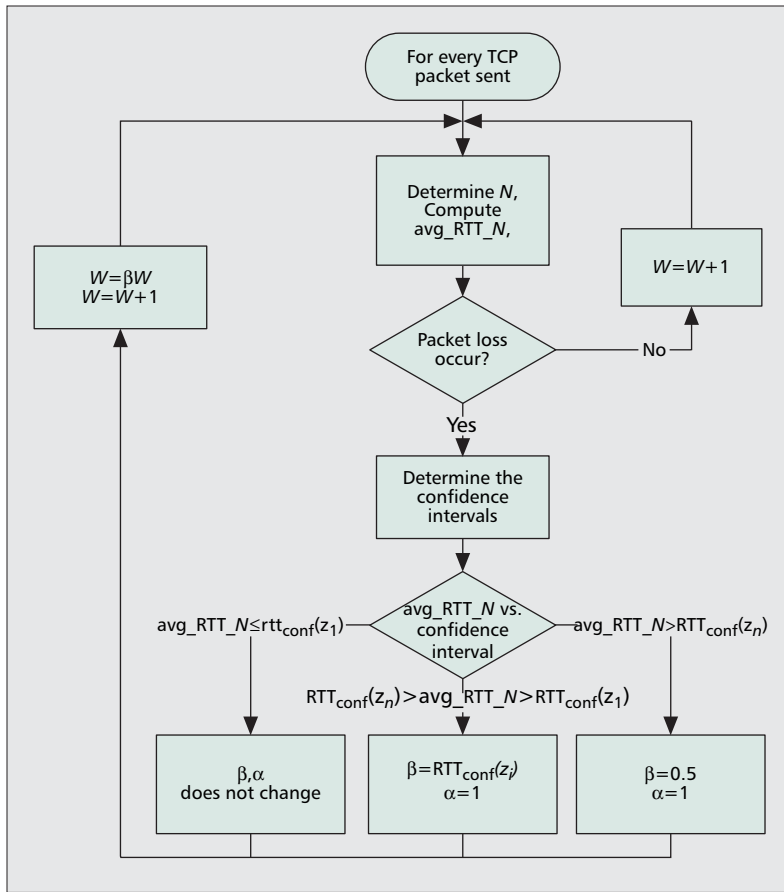
$$\alpha = \frac{4}{3}(1 - \beta^2).$$

Statistical AIMD (SAIMD) — In [40] the authors introduced a new congestion control scheme for TCP over OBS networks, called Statistical AIMD (SAIMD). SAIMD maintains and analyzes a number of previous RTTs at the TCP senders in order to identify the confidence with which a packet loss event is due to network congestion. The confidence is derived by positioning short-term RTT in the spectrum of long-term historical RTTs. The derived confidence corresponding to the packet loss is taken in the developed policy for TCP congestion window adjustment.

The SAIMD scheme adopts the framework of BAIMD to enhance the responsiveness of TCP to any burst loss event that is not caused by congestion. In SAIMD, when a burst consisting of many TCP segments from single or multiple TCP senders is lost, the corresponding TCP senders are notified of the packet loss through receiving either TDed *acks* or TO. In either case, instead of halving the *cwnd* or even throttling to the slow start stage, TCP senders reduce their *cwnd* by the multiplicative factor β . The factor β is dynamically determined by positioning the short-term RTT statistics in the spectrum of long-term historical RTTs. Here, *statistics* refers to mean, standard deviation, and correlation function in this study, and are further detailed as follows.

Two parameters in the proposed scheme, denoted M and N , were introduced. M is the number of consecutive RTTs measured for long-term statistics. M should be sufficiently large that the derived statistics (i.e., the mean and standard deviation) can fully represent the intrinsic characteristics of the network topology, routing policy, and traffic distribution/pattern. N is the number of consecutive RTTs measured prior to a packet loss for the short-term statistics. The average of N RTTs, denoted avg_RTT_N , is compared with the average of M RTTs, denoted avg_RTT_M , in a TCP session in order to determine how likely that the packet loss is due to network congestion or random burst contention loss in a lightly loaded OBS network. In a packet loss event caused by random burst contention, avg_RTT_N is expected to be close to avg_RTT_M since the high utilization of network resources remains only a short time period in the N RTTs. A larger avg_RTT_N can indicate that a packet loss event is more likely due to network congestion rather than random burst contention.

The relationship between avg_RTT_M and avg_RTT_N is based on the following observations:



■ Figure 7. Flowchart of the SAIMD congestion control scheme.

- In TCP over OBS networks, packet losses can be caused by random burst contention in OBS core networks or network congestion along the route of IP access networks and OBS core networks. The difference between random burst contention and network congestion is that network congestion suffers from high resource utilization for a longer period.
- In the high resource utilization state, the RTT of each packet delivery will be much higher than that in the low utilization state. This is due to the fact that high utilization causes longer queuing delay in IP access networks. Also, in an OBS core network with contention resolution schemes such as burst retransmission [24, 25] and deflection [26, 27], bursts have a high probability of being retransmitted or deflected, which results in a longer average burst delay in the OBS network.

Since N RTTs are expected to provide sufficient information about the short-term network status when a packet is lost, the scheme uses autocorrelation to select a proper value of N . If N is chosen too small or too large, the short-term network status may not be accurately represented. The scheme defines an autocorrelation threshold value (γ) that determines N . In order to represent the short-term network status well, the N RTTs should have strong correlation with each other. Hence, the value of γ should be close to 1. In SAIMD γ is taken as 90 percent. Figure 7 shows the congestion control of the proposed SAIMD scheme.

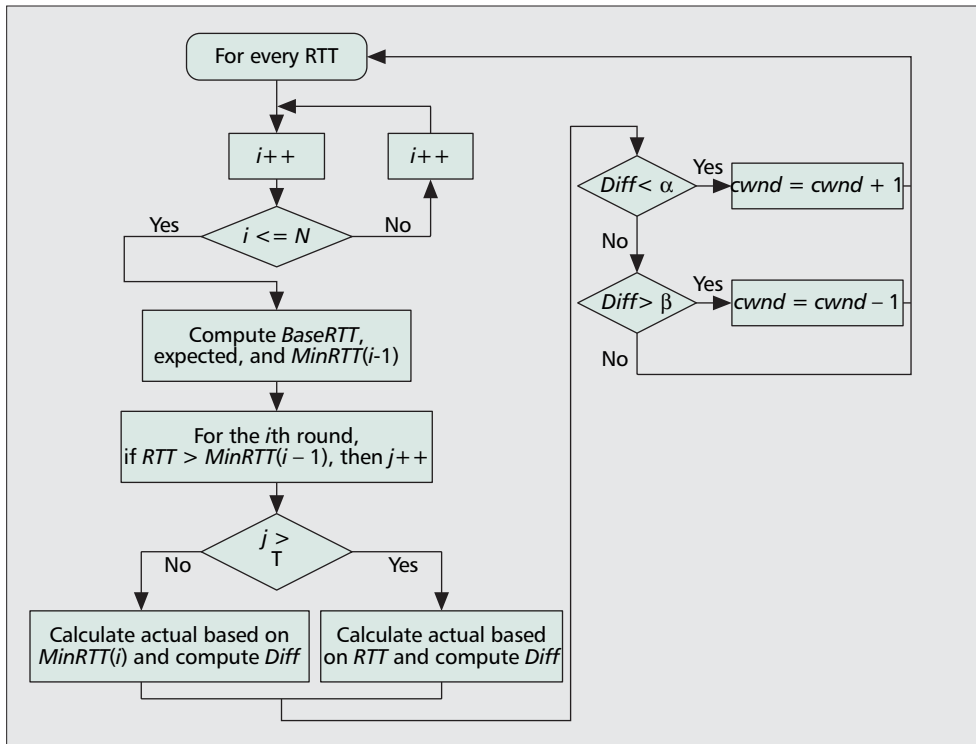
Compared to the conventional AIMD-based TCP scheme, SAIMD causes additional overhead for maintaining the M RTTs along with the efforts in computing the autocorrelation and confidence intervals for the N RTTs. The cost is nonetheless a trade-off with the long convergence time in recovery from slow start caused by false congestion detection. This is

considered of essential importance for those high-bandwidth TCP flows that may take hours or days to recover from a slow start. Note that the computation for the autocorrelation and confidence interval is required only when a segment loss event occurs, and the computation complexity is almost a constant regardless of M and N . In addition, the proposed SAIMD scheme can mainly be used for long and high-bandwidth TCP flows instead of short TCP such as HTTP Web services; thus, the resultant additional overhead in the whole network is expected to be trivial.

Threshold-Based TCP Vegas — In [36, 37] the authors analyzed the performance of delay-based TCP (i.e., Vegas) over OBS with burst retransmission and deflection. The study shows several issues observed when the conventional TCP Vegas congestion control mechanism is adopted in OBS networks. When a fixed source routing strategy is used, the packet delay experienced in the OBS domain is primarily the sum of burst assembly delay and link propagation delay, which do not vary when the traffic load changes in OBS networks. Hence, conventional TCP Vegas cannot effectively detect network congestion in the OBS domain. Furthermore, if all packets in the congestion window of TCP Vegas are assembled into a single burst, TCP Vegas may suffer from the false congestion detection problem, which fatally impairs TCP throughput due to unnecessary TO retransmissions followed by slow start at TCP Vegas senders.

When TCP Vegas runs over OBS networks with burst retransmission or deflection, the false congestion detection problem can be mitigated since both schemes incur extra delay for retransmitted bursts. Hence, TCP Vegas will detect the increases in RTTs for packets in bursts that are retransmitted, which may result in TCP Vegas reducing $cwnd$, leading to lower TCP throughput. If the increases in RTTs are caused by burst retransmission in a lightly loaded OBS network, TCP Vegas should not reduce its $cwnd$. Hence, a modified TCP Vegas able to tell whether the increase in RTTs is due to network congestion or to retransmission in lightly loaded OBS networks is introduced and evaluated in [36, 37].

Based on the above observations, the authors proposed a threshold-based TCP Vegas scheme by introducing a new parameter, T , referred to as the threshold, to assist TCP Vegas to distinguish between network congestion and burst contention at low traffic loads. In the proposed scheme TCP Vegas measures RTT for each packet sent and keeps track of the minimum measured RTTs of N consecutive packets. Let $MinRTT(i)$ be the minimum measured RTTs of i ($0 < i < N$) consecutive packets. In the i th round, if the measured RTT of the i th packet is larger than $MinRTT(i-1)$, it means that the i th packet was once queued in the access network and/or assembled in a burst that was retransmitted. A counter that keeps the number of packets whose RTTs are larger than their $MinRTT(i-1)$ will then be increased by 1. If the number of TCP packets whose RTTs are larger than their $MinRTT(i-1)$ is under the threshold T , TCP sender will stay with the *Actual* throughput calculated based on $MinRTT(i-1)$, even if the measured RTTs are increased. If the number of TCP packets whose RTTs are larger than their $MinRTT(i-1)$ exceeds the threshold T , it means that the network is congest-



■ **Figure 8.** The threshold-based TCP Vegas congestion control scheme.

ed. Hence, threshold-based Vegas recognizes the network congestion and calculates the *Actual* throughput as usual.

The fundamental object of threshold-based TCP Vegas is to reduce the sensitivity of TCP Vegas to increased RTTs caused by burst retransmission in the OBS domain. Instead of changing $cwnd$, the sensitivity of TCP Vegas can also be reduced by decreasing α or increasing β . However, changing α and β makes it difficult for TCP Vegas to estimate the available bandwidth in the networks. Figure 8 summarizes the proposed threshold-based TCP Vegas.

In threshold-based TCP Vegas, the number of packets consecutively sent (denoted N) and threshold T should be chosen much larger than the number of packets from a single TCP Vegas connection that are assembled into a burst so that TCP Vegas is able to detect the frequency of retransmission in the OBS domain based on a number of bursts. Usually, the packets from a TCP connection assembled in the same burst have the same measured RTT. By analyzing the variation pattern of packet RTTs, TCP Vegas can obtain the number of packets from a TCP Vegas connection that are assembled into a burst. Also, the relationship between the T and N affects the TCP performance. When T is chosen closer to N , there would be fewer remaining packets in the N consecutive packets to react to the detected congestion, which results in an ineffective response to network congestion. Hence, we take $N = iT$, where $i > 1$. Defining proper values of T and N is subject to many factors such as the number of TCP flows, TCP flow speed, and burst assembly thresholds. Optimal values of T and N are derived in [50] using the fixed-point feedback mechanism proposed in [51].

Table 2 summarizes the features of the surveyed schemes in this review article; note that BCR stands for solutions that work on OBS with burst contention resolution schemes.

OPEN RESEARCH PROBLEMS

Although TCP has been subject to extensive research efforts in the past decades, TCP over OBS is a relatively unexplored

research area with a limited number of studies that tackled some of the unique features in such a network scenario. Through close analysis of the reported TCP enhancements listed in Tables 1 and 2, we observed that the most important characteristics and abilities for a congestion control mechanism in TCP over OBS lie in the following:

- Able to handle multiple TCP segment loss events in one round-trip
- Able to identify false TOs and burst losses under low traffic load situation
- Able to compensate for out-of-order delivery in the OBS domain

In this section we identify the open issues in the area of TCP over OBS.

INTEGRATING LINK LAYER SOLUTIONS WITH TCP PERFORMANCE

Although there has been less effort on TCP enhancements over OBS, there have been extensive studies reported on the OBS network that aim to reduce burst dropping, provide QoS in burst transmission, and conduct cross-layer design optimization in terms of burst assembly delay, burst size, burst delivery, and so on. The link layer schemes, such as adaptive assembly algorithms [14], burst retransmission [24, 25], burst deflection [26, 27], and FDLs [15], have successfully contributed in reducing burst loss probability at the expense of introducing extra delay and design complexity. It is important to integrate link layer solutions with any mechanism for TCP throughput and reliability enhancements before OBS can be practically deployed in the Internet.

We found that reducing the burst dropping probability may not result in all cases in significant enhancement in TCP throughput. For example, burst retransmission in OBS networks can greatly improve TCP performance; however, the persistence of retransmission, deflection, or segmentation is subject to further considerations since too much persistence leads to TO in the TCP layer, which significantly throttles TCP transmission. One of the most important

Schemes	Solution category	OBS type	Devices involved			Explicit notifications	Problem addressed		
			TCP source	OBS edge	OBS core		Random burst loss	Slow convergence	Packet reordering
Adaptive Assembly Period (AAP)	Link layer	Barebone/BCR		√				√	√
Burst Retransmission	Link layer	BCR		√	√		√		
Deflection Routing	Link layer	BCR			√		√		
FDLs	Link layer	BCR			√		√		
Burst Segmentation	Link layer	BCR			√		√		
Wavelength Conversions	Link layer	BCR			√		√		
Retransmission-Count Dropping	Link layer	BCR		√	√		√	√	
TCP with Burst Acknowledgment	Link layer	Barebone/BCR	√	√	√	√	√		
TCP Decoupling	Link layer	Barebone/BCR	√	√			√		
BTCP with Burst Length Estimation	Without signaling	Barebone/BCR	√						
BTCP with burst ack	Signaling	Barebone/BCR	√	√		√	√		
BTCP with burst Nack	Signaling	Barebone/BCR	√	√	√	√	√		
Burst AIMD	Without signaling	Barebone/BCR	√	√	√		√	√	
TCP-Burst Contention Loss	Signaling	Barebone/BCR	√	√	√	√	√	√	√
Statistical AIMD	Without signaling	Barebone/BCR	√				√	√	√
Threshold-based Vegas	Without signaling	BCR	√				√		√

■ Table 2. Overview of the surveyed OBS solutions for TCP.

challenges while tackling this problem is to properly define the TO threshold in the TCP senders in the presence of various traffic loads at edge and core nodes. Determination of the threshold value should take into consideration the number of successively assembled packets in a single burst. A TCP snooping mechanism similar to some proposed schemes in wireless networks, such as ATCP [52], ELN [10], JTCP [53], TCP Casablanca [54], and TCP Peach [55], can be developed to evaluate the persistence of burst retransmission by estimating the RTT. This enhancement aims to compromise the risk of having a false TO with the gain through performing burst retransmission or burst deflection routing.

The research on TCP over OBS presented in BTCP [32], BAIMD [39], TCP-BCL [35], and SAIMD [40] has successfully resolved the vicious effects of false congestion detection due to the bufferless characteristic of OBS networks. They improved the TCP sender reaction to each packet loss event relevantly. However, each scheme suffers from some overhead as well, such as high computation complexity and/or extra signaling in the OBS networks. Also, some of the above schemes may fail to overcome the problem of losing a large number of packets assembled in a single burst. Furthermore, it is necessary to evaluate the convergence rate of TCP in the presence of either a constant RTT in the barebone OBS or different values of RTT due to the deployment of burst contention resolution schemes.

TCP CONVERGENCE RATES FOR LARGE BANDWIDTH-DELAY PRODUCT NETWORKS

There is a significant lack of performance evaluation on the TCP modifications proposed for a network environment with a large bandwidth-delay product over OBS networks. Fast TCP [7, 8], binary increase congestion control (BIC) TCP [56], Explicit Control Protocol (XCP) [57], TCP with Simple Available Bandwidth Utilization Library (SABUL) [58], High Speed (HSTCP) [59], and Scalable TCP (STCP) [60] are among the most famous promising solutions. XCP has shown stability and efficiency using ECN through extending ECN and Core Stateless Fair Queue (CSFQ), which make XCP aware of the per-flow state and buffer size status. XCP uses multiplicative increase multiplicative decrease (MIMD) to control the congestion window, which increases the transmission by Δ (i.e., $\Delta =$ the bandwidth squared – the queue size). On the other hand, it uses AIMD to control the fairness. If $\Delta > 0$, divide by Δ equally between the coexisting flows; otherwise, divide Δ between flows proportionally to their current rates [57]. SABUL introduces a hybrid approach by merging the rate-based transmission via UDP and reliable retransmission via TCP, where a UDP channel is adopted for transmitting data at high rates, while a TCP channel is used to resend the missing data segments to ensure reliability [55]. Depending on the current *cwnd*, HSTCP uses $a(cwnd)$ and $b(cwnd)$ for computing the next window size. This scheme is known to be a safe incremental approach [59]. A simulation-based study for HSTCP over OBS is presented in [33]. Using small burst assembly delay and moderate burst dropping, the study shows that HSTCP throughput is severally affected. Scalable TCP (STCP) uses MIMD approach with sending rate 0.01 and multiplicative rate as 1/8. STCP is a sender-side TCP that offers a robust mechanism to improve performance in high-speed wide area networks using traditional TCP receivers. STCP increases the TCP's *cwnd* by 0.01 for each acknowledged packet (not in the fast recovery stage) and cuts the *cwnd* by 0.875 for each packet loss event [60].

There are several proposed schemes that solve the TCP slow convergence problem over high-speed networks. There is a need to evaluate the burst assembly delay and burst dropping over these TCP congestion control algorithms. There are great opportunities for investigating the effect of the burst assembly delay, burst dropping vs. packet aggregation gain on Fast TCP, Scalable TCP, XCP, and SABUL. It is known that the above schemes can achieve faster convergence of TCP throughput in large-bandwidth high-delay networks by quickly enlarging their *cwnd*. However, with large *cwnd*, the number of *acks* is significantly decreased. The presence of random burst losses that contain a large number of *ack* packets results in dramatic damage to TCP *ack*-clocking and forces TCP to fall into false detection of network congestion, thus retransmitting an unnecessarily large number of packets. Furthermore, the *ack* losses are expected to affect large number of TCP senders since the burst can assemble many *ack* packets due to their small size. Up to our knowledge, the effect of *ack* packet losses over OBS networks has not been addressed in the literature.

PERFORMANCE MODELING FOR TCP OVER OBS NETWORKS

The previously reported TCP over OBS performance modeling technique follows the packet-oriented approach [24, 30, 31, 35, 37, 40]. These studies use Markov modulated Poisson arrival of the traffic [61]. In the early 1990s, the fluid modeling technique proposed in [62] added a new dimension to modeling large numbers of TCP flows. However, the fluid modeling approach requires maintaining a few strict assumptions, such as:

- Having a very large number of TCP flows
- Poisson arrival of loss events
- Strong correlation between losses in one RTT while being independent among the other RTTs

Regarding the first assumption, there is no sufficient evidence that the number of TCP flows is sufficiently high at the OBS edge node. In OBS, since both random burst drops and dropping due to persistent congestion may occur, the second assumption is subject to further investigation. The third assumption can partially be justified as in the barebone OBS the RTT is more or less fixed. This is because the third assumption can only hold for TCP flows that can emit the TCP segments of their entire *cwnd* while being assembled in one burst (e.g., fast flows [30]). We conclude that using the fluid model for evaluating TCP throughput performance over OBS requires significant improvements before it is considered accurate.

The synchronization modeling approach proposed in [63] benefits from *ack*-clocking to include the burstiness factor in the fluid model. Note that the fluid model assumes that there is no burstiness, and the TCP rates of different flows are differentiable. Therefore, it may take an infinitely long time to converge. The synchronization approach has been used for modeling Fast TCP and obtaining its stability by [63, 64]. In order to obtain sufficient analysis for TCP performance while considering TCP stability, scalability, and responsiveness, the synchronization modeling approach needs to capture the bufferless nature of OBS links (i.e., fixed TCP RTTs), the burst aggregation factors, and the burst loss distribution.

CONCLUSIONS

The research efforts in modifying and extending conventional TCP implementations have been extensively reported in past years; however, relatively limited knowledge has been gained in terms of the impacts on the upper layer applications when OBS is adopted underneath in the presence of burst dropping at different traffic loads. This article provides a comprehensive review and tutorial on the previously reported congestion control schemes for TCP over OBS networks. With the taxonomy identified in this review, these schemes are classified into three categories, including link layer solutions, and congestion detection mechanisms with and without explicit notifications between the OBS edge nodes and TCP senders. We have enumerated and discussed each category in detail, aiming to provide a complete picture of the development of TCP enhancements over OBS networks. We have also presented some future research directions for TCP enhancements in the presence of various OBS characteristics, including the scenarios of large bandwidth-delay products, schemes for protecting TCP *ack* packets, and extensions and elaborations of the existing modeling techniques (e.g., fluid and synchronization modeling), in order to capture the burst transmission characteristics of TCP over OBS networks.

REFERENCES

- [1] W. Stevens, *TCP/IP Illustrated, Volume 1*, Addison Wesley Longman, 1994.
- [2] J. Postel, "Transmission Control Protocol," RFC 793, Protocol Spec., DARPA Internet Program, 1981
- [3] C. Fraleigh et al., "Packet-Level Traffic Measurements from the Sprint IP Backbone," *IEEE Network*, vol. 17, no. 6, Nov.-Dec., 2003, pp. 6–16.
- [4] W. Stevens, "TCP Slow Start, Congestion Avoidance, Fast Retransmit, and Fast Recovery Algorithms," RFC 2001, 1997.
- [5] M. Mathis et al., "TCP Selective Acknowledgment Options," RFC 2018, 1996.

- [6] L. Brakmo and L. Peterson, "TCP Vegas: End-to-End Congestion Avoidance on A Global Internet," *IEEE JSAC*, vol. 13, no. 8, Oct. 1995, pp. 1465–80.
- [7] D. Katabi, M. Handley, and C. Rohrs, "Congestion Control for High Bandwidth-Delay Product Networks," *Proc. ACM SIGCOMM*, 2002.
- [8] C. Jin, D. Wei, and S. Low, "FAST TCP: Motivation, Architecture, Algorithms, Performance," *Proc. IEEE INFOCOM*, 2004.
- [9] S. Hegde et al., "FAST TCP in High-Speed Networks: An Experimental Study," *Wksp. Networks for Grid Apps.*, Oct. 2004.
- [10] H. Balakrishnan and R. Katz, "Explicit Loss Notification and Wireless Web Performance," *Proc. IEEE GLOBECOM*, Internet MiniConf., 1998.
- [11] S. Floyd, "TCP and Explicit Congestion Notification," *ACM SIGCOMM Comp. Commun. Rev.*, vol. 24, no. 5, 1994, pp. 8–23.
- [12] T. Battestilli and H. Perros, "An Introduction to Optical Burst Switching," *IEEE Optical Commun.*, vol. 41, 2003, pp. 510–15.
- [13] Y. Xiong, M. Vandenhoue, and H. Cankaya, "Control Architecture in Optical Burst-Switched WDM Networks," *IEEE JSAC*, vol. 18, no. 10, 2000, pp. 1838–51.
- [14] X. Cao et al., "Assembling TCP/IP Packets in Optical Burst Switched Networks," *Proc. IEEE GLOBECOM*, 2002.
- [15] I. Chlamtac et al., "CORD: Contention Resolution by Delay Lines," *IEEE JSAC*, vol. 14, no. 5, 1996, pp. 1014–29.
- [16] Y. Chen, C. Qiao, and X. Yu, "An Optical Burst Switching: A New Area in Optical Networking Research," *IEEE Network*, vol. 18, no. 5, 2004, pp. 16–23.
- [17] M. Yoo and C. Qiao, "Just-Enough-Time (JET): A High-Speed Protocol for Bursty Traffic in Optical Networks," *Proc. IEEE/LEOS Tech. for Global Info. Infrastructure*, 1997, pp. 26–27.
- [18] J. Y. Wei and R. I. McFarland, "Just-In-Time Signaling for WDM Optical Burst Switching Networks," *J. Lightwave Tech.*, vol. 18, 2000, pp. 2019–37.
- [19] J. White, R. Tucker, and K. Long, "Merit-based Scheduling Algorithm for Optical Burst Switching," *Int'l. Conf. Optical Networks*, Korea, 2002, pp. 75–77.
- [20] I. Ogushi et al., "Parallel Reservation Protocols for Achieving Fairness in Optical Burst Switching," *Proc. IEEE Wksp. High Perf. Switching and Routing*, 2001, pp. 213–17.
- [21] A. Kaheel and H. Alnuweiri, "Batch Scheduling Algorithms for Optical Burst Switching Networks," *Proc. IFIP Networking*, 2005, pp. 90–101.
- [22] B. Zhou, M. Bassiouni, and G. Li, "Improving Fairness in Optical-Burst-Switching Networks," *J. Optical Networking*, vol. 3, no. 4, 2004, p. 214.
- [23] J. Li et al., "Maximizing Throughput for optical Burst Switching Networks," *Proc. IEEE INFOCOM*, 2004, pp. 1853–63.
- [24] Q. Zhang et al., "Evaluation of Burst Retransmission in Optical Burst-Switched Networks," *Proc. 2nd Int'l. Conf. Broadband Networks '05*, Boston, MA, 2005.
- [25] Q. Zhang et al., "Analysis of TCP over Optical Burst-Switched Networks with Burst Retransmission," *Proc. IEEE GLOBECOM*, St. Louis, MO, 2005.
- [26] C. Hsu, T. Liu, and N. Huang, "Performance Analysis of Deflection Routing in Optical Burst-Switched Networks," *Proc. IEEE INFOCOM*, vol. 1, 2002, pp. 66–73.
- [27] M. Schlosser, E. Patzak, and P. Gelpke, "Impact of Deflection Routing on TCP Performance in Optical Burst Switching Networks," *7th Int'l. Conf. Transparent Optical Networks*, 2005.
- [28] V. Vokkarane, J. Jue, and S. Sitarman, "Burst Segmentation: An Approach for Reducing Packet Loss in Optical Burst Switched Networks," *IEEE ICC*, vol. 5, 2002, pp. 2673–77.
- [29] A. Zalesky et al., "Evaluation of Limited Wavelength Conversion and Deflection Routing as Methods to Reduce Blocking Probability in Optical Burst Switched Networks," *IEEE ICC*, vol. 3, 2004, pp. 1543–47.
- [30] A. Detti and M. Listanti, "Impact of Segment Aggregation on TCP Reno Flows in Optical Burst Switching Networks," *Proc. IEEE INFOCOM*, 2002.
- [31] X. Yu et al., "Performance Evaluation of TCP Implementations in OBS Networks," tech. rep. 2003-13, SUNY Buffalo, 2003.
- [32] X. Yu, C. Qiao, and Y. Liu, "TCP Implementation and False Time Out Detection in OBS Networks," *Proc. IEEE INFOCOM*, 2004.
- [33] L. Zhu, N. Ansari, and J. Liu, "Throughput of High-Speed TCP in Optical Burst Switching Networks," *IEE Proc. Commun.*, vol. 152, no. 3, 2005, pp. 349–52.
- [34] P. Du and S. Abe, "TCP Performance Analysis of Optical Burst Switching Networks with a Burst Acknowledgment Mechanism," *Asia-Pacific Conf. Commun. and Int'l. Symp. Multi-Dimensional Mobile Commun.*, Aug. 2004.
- [35] B. Shihada and P.-H. Ho, "A Novel TCP with Dynamic Burst-Contention Loss Notification over OBS Networks," *Elsevier J. Comp. Networks*, vol. 52/2, 2008, pp. 461–71.
- [36] B. Shihada, Q. Zhang, and P.-H. Ho, "Threshold-based TCP Vegas over Optical Burst Switched Networks," *15th IEEE Int'l. Conf. Comp. Commun. and Networks*, VA, Oct. 2006.
- [37] B. Shihada, Q. Zhang, and P.-H. Ho, "Performance Evaluation of TCP Vegas over Optical Burst Switched Networks," *IEEE Broadnets*, San Jose, CA, Oct. 2006.
- [38] S. Gowda et al., "Performance Evaluation of TCP over Optical Burst-Switched (OBS) WDM Networks," *Proc. IEEE ICC*, 2003, pp. 1433–37.
- [39] B. Shihada et al., "BAIMD: a Responsive Rate Control for TCP over Optical Burst Switched (OBS) Networks," *IEEE ICC*, Turkey, 2006.
- [40] B. Shihada, P.-H. Ho, and Q. Zhang, "SAIMD: A Congestion Detection Scheme for TCP over OBS Networks," accepted in *IEEE J. Lightwave Tech.*, Oct. 2007.
- [41] S. Y. Wang, "Using TCP Congestion Control to Improve the Performances of Optical Burst Switched Networks," *Proc. IEEE ICC*, 2003.
- [42] L. Kim, S. Lee, and J. Song, "Dropping Policy for Improving the Throughput of TCP over Optical Burst-Switched Networks," *ICOIN*, 2006, pp. 409–18.
- [43] S. Lee and L. Kim, "Drop Policy to Enhance TCP Performance in OBS Networks," *IEEE Commun. Letters*, vol. 10, no. 4, 2006.
- [44] S. Bhandarkar and N. Reddy, "Improving the Robustness of TCP to Non-Congestion Events," Internet draft, draft-ietf-tcpmp-tcp-dcr-01.txt, 2004.
- [45] P. Sarolahti and M. Kojo, "F-RTO: an Algorithm for Detecting Spurious Retransmission Timeouts with TCP and SCTP," Internet draft, draft-ietf-tcpmp-frto-01.txt, 2004.
- [46] R. Ludwig, A. Gurtov, "The EIFEL Response Algorithm for TCP," Internet draft, draft-ietf-tsvwg-tcp-eifel-response-05.txt, 2004.
- [47] S. Bhandarkar et al., "TCP-DCR: A Novel Protocol for Tolerating Wireless Channel Errors," *IEEE Trans. Mobile Comp.*, 2004.
- [48] Y. Yang and S. Lam, "Generalized AIMD Congestion Control," Univ. TX, tech. rep. TR-2000; <http://www.cs.utexas.edu/users/lam/NRL/TechReports/>
- [49] L. Cai et al., "Performance Analysis of TCP-Friendly AIMD Algorithms for Multimedia Applications," *IEEE Trans. Multimedia*, vol. 7, no. 2, 2005, pp. 339–55.
- [50] B. Shihada, Q. Zhang, and P.-H. Ho, "Performance Evaluation of Threshold-Based TCP Vegas over Optical Burst Switched Networks," before *IEEE JSAC*, Jan. 2007.
- [51] C. Cameron et al., "TCP over OBS — Fixed-Point Load and Loss," *Optics Express*, vol. 13, 2005, pp. 9167–74.
- [52] J. Liu, and S. Singh, "ATCP: TCP for Mobile Ad Hoc Networks," *IEEE JSAC*, vol. 19, no. 9, 2001, pp 1300–15.
- [53] E. H.-K. Wu and M.-Z. Chen, "JTCP: Jitter-based TCP for Heterogeneous Wireless Networks," *IEEE JSAC*, vol. 13, no. 4, 2004, pp. 757–66.
- [54] S. Biaz and N. Vaidya, "De-Randomizing Congestion Losses to Improve TCP Performance over Wired-Wireless Networks," *IEEE/ACM Trans. Net.*, vol. 13, no. 3, 2005, pp. 596–608.
- [55] I. Akyildiz, G. Morabito, and S. Palazzo, "TCP-Peach: A New Congestion Control Scheme for Satellite IP Networks," *IEEE/ACM Trans. Net.*, vol. 9, no. 3, 2001, pp. 307–21.
- [56] L. Xu, K. Harfoush, and I. Rhee, "Binary Increase Congestion Control (BIC) for Fast Long-Distance Networks," *IEEE INFOCOM*, Hong Kong, China, 2004.
- [57] A. Falk and D. Katabi, "XCP Protocol Specification," Internet draft, Feb. 2004.
- [58] H. Sivakumar et al., "Simple Available Bandwidth Utilization Library for High-Speed Wide Area Networks," *J. Supercomputing*, 2004.
- [59] Floyd, "High-Speed TCP for Large Congestion Windows," RFC 3649, experimental, Dec. 2003.
- [60] T. Kelly, "Scalable TCP: Improving Performance in High-Speed Wide Area Networks," *ACM SIGCOMM*, vol. 33, no. 2, 2003,

-
- pp. 83–91.
- [61] T. Karagiannis *et al.*, “A Nonstationary Poisson View of Internet Traffic,” *IEEE INFOCOM*, 2004.
- [62] V. Misra, W.-B. Gong, and D. Towsley, “Fluid Based Analysis of a Network AQM Routes Supporting TCP Flows with an Application to RED,” *Proc. ACM SIGCOMM*, 2000, pp. 151–60.
- [63] D. Wei, “Congestion Control Algorithms for High-Speed Long Distance TCP Connections,” Master thesis, Caltech, 2004; <http://www.cs.caltech.edu/~weixl/research/msthesis.ps>
- [64] J. Wang, D. Wei, and S. Low, “Modeling and Stability of Fast TCP,” *IEEE INFOCOM*, 2005, pp. 938–48.

BIOGRAPHIES

BASEM SHIHADA (bshihada@uwaterloo.ca) received his B.Sc., M.Sc., and Ph.D. in computer science from United Arab Emirates University in 1997, Dalhousie University in 2001, and the University of Waterloo in 2007, respectively. His research interests are in high-speed optical networks, optical burst switching, grid computing,

transport layer protocol, resource management, and system security.

PIN-HAN HO (pinhan@bbcr.uwaterloo.ca) received his B.Sc. and M.Sc. degrees from the Electrical and Computer Engineering Department of National Taiwan University in 1993 and 1995. He started his Ph.D. study in 2000 at Queen’s University, Kingston, Canada, focusing on optical communications systems, survivable networking, and QoS routing problems. He finished his Ph.D. in 2002, and joined the Electrical and Computer Engineering (ECE) Department of the University of Waterloo, Canada, as an assistant professor the same year. He is the author/coauthor of more than 100 refereed technical papers and book chapters, and co-author of a book on optical networking and survivability. He is the recipient of the Distinguished Research Excellence Award from the ECE Department of the University of Waterloo, the Early Researcher Award (Premier Research Excellence Award) in 2005, Best Paper Awards at SPECTS ’02, ICC ’05 Optical Networking Symposium, and ICC ’07 Security and Wireless Communications Symposium, and the Outstanding Paper Award at HPSR ’02.