

# On TCP Performance in a Heterogeneous Network: A Survey

Chadi Barakat, Eitan Altman, and Walid Dabbous, INRIA

## ABSTRACT

Transmission media carrying Internet traffic present a wide range of characteristics, some of which, such as transmission errors, long end-to-end delay, and bandwidth asymmetry, may cause a degradation of TCP performance. Many works have studied the performance of TCP over these media, most of which focus on a particular network type. In this work we study TCP performance independent of the type of network by considering the different possible characteristics of the connection path. We present the problems and the different proposed solutions. This study permits us to understand the limitations of the actual solutions and the required modifications to let TCP cope with a heterogeneous Internet on an end-to-end basis.

## INTRODUCTION

The Transmission Control Protocol (TCP) provides a reliable connection-oriented in-order service to many of today's Internet applications. Given the simple best-effort service provided by IP, TCP must cope with the different transmission media crossed by Internet traffic. This mission of TCP is becoming difficult with the increasing heterogeneity of the Internet. High-speed links (optic fibers), long and variable delay paths (satellite links), lossy links (wireless networks), asymmetric paths (hybrid satellite networks), and others are becoming widely embedded in the Internet. Many works have studied — by experimentation [1–4], analytical modeling [5–8], and simulation [6, 9–11] — the performance of TCP in this new environment. Most of these works have focused on a particular environment (satellite networks, mobile networks, etc.). They have revealed some problems in the operation of TCP. Long propagation delay and losses on a satellite link, handover and fading in a wireless network, bandwidth asymmetry in some media, and other phenomena have been shown to seriously affect the throughput of a TCP connection. A large number of solutions have been proposed. Some solutions suggest modifications to TCP to help it to cope with these new paths. Other solutions keep the protocol unchanged and hide the problem from TCP. Typically, modifications which

vary according to the problem to hide are proposed inside the network.

In this article we summarize the different works on TCP performance. Our aim is to conduct this study independent of the type of network causing the problem. Instead of talking about particular environments, we consider the different characteristics of a path crossed by TCP traffic, focusing on bandwidth-delay product (BDP), round-trip time (RTT), noncongestion losses, and bandwidth asymmetry. After overviewing TCP, we study the impact of each of these characteristics on TCP, and present the proposed solutions as well as our comments.

## AN OVERVIEW OF TCP

TCP is a reliable window-based acknowledgment (ACK)-clocked flow control protocol. It uses an additive-increase multiplicative-decrease strategy for changing its window as a function of network conditions [4, 12]. Starting from one packet, or a larger value as we will see later, the window is increased exponentially by one packet for every nonduplicate ACK until the source estimate of network capacity is reached. By capacity of the network, sometimes called the pipe size, we mean the maximum number of packets that can be fit on the path. This is the slow-start (SS) phase, and the capacity estimate is called the SS threshold ( $ss_{thresh}$ ). SS aims to alleviate the burstiness of TCP while quickly filling the pipe. Once  $ss_{thresh}$  is reached, the source switches to a slower increase in the window by one packet for every window's worth of ACKs. This phase, called congestion avoidance (CA), aims to slowly probe the network for any extra bandwidth. The window increase is interrupted when a loss is detected. Two mechanisms are available for the detection of losses: the expiration of a retransmission timer (timeout) or the receipt of three duplicate ACKs (fast retransmit, FRXT). The source supposes that the network is in congestion and sets its estimate of the capacity to half the current window. Tahoe [4, 9], the first version of TCP to implement congestion control, at this point sets the window to one packet and uses SS to reach the new  $ss_{thresh}$ . Slow starting after every loss detection deteriorates the performance given the low bandwidth utilization during SS. When

快速重传也要窗口减半？

the loss is detected via timeout, SS is unavoidable since the ACK clock has stopped and SS is required to smoothly fill the pipe. However, in the FRXT case, ACKs still arrive at the source and losses can be recovered without SS. This is the objective of the new versions of TCP (Reno, New Reno, SACK, etc.) [9] that **call a Fast Recovery (FRCV) algorithm** to retransmit the losses while maintaining enough packets in the network to preserve the ACK clock. Once losses are recovered, this algorithm ends and normal CA is called. If FRCV fails, the ACK stream stops, a timeout occurs, and the source resorts to SS as with Tahoe.

求助于

## A LARGE BANDWIDTH-DELAY PRODUCT

The increase in link speed (e.g., optic fibers at gigabits per second) has led to paths of large BDP. TCP window must be able to reach large values in order to efficiently use the available bandwidth. Large windows, up to  $2^{30}$  bytes, are now possible with the window scale TCP option [13]. However, at large windows, congestion may lead to the loss of many packets from the same connection. Efficient FRCV is then required to correct many losses from the same window. Also, at large BDP, network buffers have an important impact on performance. These buffers must be well dimensioned and scale with the BDP.

### FAST RECOVERY

FRCV uses the information carried by ACKs to estimate the number of packets in flight while recovering from losses. New packets are sent if this number falls below the network capacity estimate. The objective is to preserve the ACK clock in order to avoid the timeout.

The difference between the different versions of TCP **is in the estimation of the number of packets in flight during FRCV**. Reno [9] considers every duplicate ACK a signal that a packet has left the network. The problem of Reno is that it leaves FRCV when an ACK for the first loss in a window is received. This prohibits the source from detecting the other losses with FRXT [9]. A long timeout is required to detect the other losses. New Reno [9, 11] has been proposed to overcome this problem. The idea is to stay in FRCV until all the losses in the same window are recovered. Partial ACKs are used to detect multiple losses in the same window. This avoids timeout but cannot result in a recovery faster than one loss per RTT. The source needs to wait for the ACK of the retransmission to discover the next loss. Another problem of Reno and New Reno is that they rely on ACKs to estimate the number of packets in flight. ACKs can be lost on the return path, which results in an underestimation of the number of packets that have left the network, and thus an underutilization of the bandwidth during FRCV and, in the case of Reno, a possible failure of FRCV.

More information is needed at the source to recover faster than one loss per RTT and to estimate more precisely the number of packets in the pipe. This information is provided by

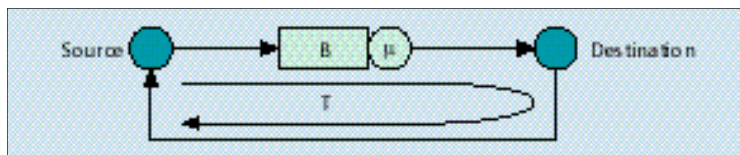


Figure 1. The single-node network model.

selective ACK (SACK) [13], a TCP option containing the three blocks of contiguous data most recently received at the destination. Many algorithms have been proposed to use this information during FRCV. We find TCP-SACK [9], which uses ACKs to estimate the number of packets in the pipe and SACKs to retransmit more than one loss per RTT. This leads to an important improvement in performance when bursts of losses appear in the same window, but the recovery is always sensitive to the loss of ACKs. As a solution, we find forward ACK (FACK) [12], which relies on SACK in estimating the number of packets in the pipe. The number and identity of packets to transmit during FRCV is decoupled from the ACK clock, in contrast to TCP-SACK, where the identity is only decoupled.

### THE IMPACT OF BUFFER SIZE ON TCP PERFORMANCE

SS results in bursts of packets sent at a rate exceeding the bottleneck bandwidth. When the receiver acknowledges every data packet, the rate of these bursts is equal to twice the bottleneck bandwidth. If network buffers are not well dimensioned, they will overflow early during SS before reaching the network capacity. This will result in an underestimation of the available bandwidth and a deterioration in TCP performance.

Early losses during SS were first analyzed in [7]. The network is modeled with a single bottleneck node of bandwidth  $\mu$ , buffer  $B$ , and two-way propagation delay  $T$  (Fig. 1). The authors consider a long TCP-Tahoe connection where the aim of SS is to reach quickly without losses  $ssthresh$  which is equal to half the pipe size  $((B + \mu T)/2)$ . In the case of a receiver that acknowledges every data packet, they found that a buffer  $B$  larger than one third the BDP is required (larger than  $\mu T/3$ ). Their analysis can be extended to a SS phase with a different threshold, mainly to that at the beginning of the connection, where  $ssthresh$  is set to a default value. **As an example, it has been proposed in [11] to set the threshold at the beginning of the connection to the BDP** in order to switch to CA before the occurrence of losses. This will not work if the buffer is smaller than half the BDP (half  $\mu T/2$ ). In [5] we study the problem of early buffer overflow during SS for multiple routers in tandem. We show that, due to the high rate at which packets are sent during SS, queues can build up in routers preceding the bottleneck as well. Buffers in these routers must also be well dimensioned, otherwise they overflow during SS and limit the performance even though they are faster than the bottleneck. **With small buffers, losses during SS are not a signal of network con-**

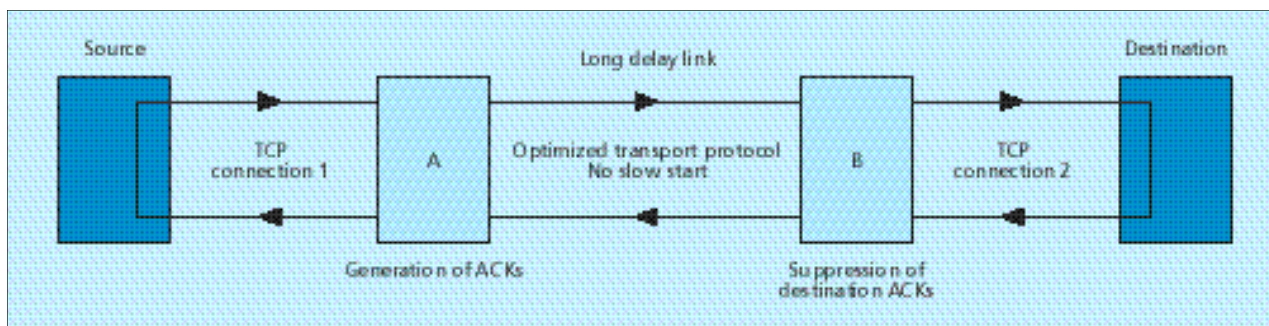


Figure 2. Spoofing: elimination of the long delay link from the feedback loop.

gestion, but rather of transient congestion due to the bursty nature of SS traffic.

Now in CA, packets are transmitted at approximately the bottleneck bandwidth. A loss occurs when the window reaches the pipe size. The source divides its window by two and starts a new cycle. To get always a throughput approximately equal to the bottleneck bandwidth, the window after reduction must be larger than the BDP. This requires a buffer B larger than the BDP. Note that we are talking about drop tail buffers, which start to drop incoming packets when the buffer is full. Active buffers such as Random Early Detection (RED) [10] start to drop packets when the average queue length exceeds some threshold. When a RED buffer is crossed by a single connection, the threshold should be larger than the BDP to get good utilization. This contrasts one of the aims of RED: limiting the size of queues in network nodes in order to reduce end-to-end delay. For multiple connections, a lower threshold is sufficient given that a small number of connections reduce their windows upon congestion, in contrast to drop tail buffers, where often all the connections reduce their windows simultaneously.

## ROUND-TRIP TIME

流行的

Long RTTs are becoming prevalent with the introduction of satellite links into the Internet. A long RTT reduces the rate at which the window increases, which is a function of the number of ACKs received and doesn't account for the RTT. This poses many problems to TCP. First, it increases the duration of SS, which is a transitory phase designed to quickly but smoothly fill the pipe [12]. Given the low bandwidth utilization during SS, this deteriorates the performance of TCP transfers, particularly short ones (e.g., Web transfers). Second, it causes unfairness in the allocation of the bottleneck bandwidth. Many works have shown the bias of TCP against connections with long RTTs [7]. Small RTT connections increase their rates more quickly and grab most of the available bandwidth. The average throughput of a connection has been shown to vary as the inverse of  $T$ , where  $\alpha$  a factor between 1 and 2 [7].

暂时的

## IMPROVING SLOW START PERFORMANCE

Many solutions have been proposed to reduce the time taken by SS on long delay paths. These solutions can be divided into three categories.

Some change the window increase algorithm of TCP, others solve the problem at the application level, and others solve it inside the network.

**TCP-Level Solutions** — The first proposition was to use a larger window than one packet at the beginning of SS [12]. An initial window of maximum four packets has been proposed. Another proposition, called byte counting, was to account for the number of bytes covered by an ACK while increasing the window rather than the number of ACKs [12]. To avoid long bursts in case of large gaps in the ACK stream, a limit on the maximum window increase has been proposed (limited byte counting).

These solutions try to solve the problem while preserving the ACK clock. They result in an increase in TCP burstiness and an overload on network buffers. Another type of solution tries to solve the problem by introducing some kind of packet spacing (e.g., rate-based spacing [12]). The source transmits directly at a large window without overloading the network. Once the large window is reached, the ACK clock takes over. This lets the source avoid a considerable part of SS.

**Application-Level Solutions** — The problem is solved at the application level without changing TCP. The first solution (e.g., XFTP [12]), consists in establishing many parallel TCP connections for the same transfer. This accelerates the growth of the resultant window, but increases the aggressiveness of the transfer and hence the losses in the network. An adaptive mechanism has been proposed for XFTP to change the number of connections as a function of network congestion.

Another solution has been proposed to accelerate the transfer of Web pages. Instead of using an independent TCP connection to fetch every object in a page, the client establishes a persistent connection and asks the server to send all the objects on it (Hypertext Transfer Protocol, HTTP, 1.1 [12]). Only the first object suffers from the long SS phase; the remaining objects are transferred at a high rate. The low throughput during SS is compensated for by the long time we stay in CA.

**Network-Level Solutions** — The problem is solved inside the network rather than at hosts, worthwhile when a long delay link is located on the path. In order to decrease the RTT, the long delay link is eliminated from the feedback loop by acknowledging packets at the input of

this link (A in Fig. 2). Packets are then transmitted on the long delay link using an optimized transport protocol (e.g., STP, proposed in [3]). This transport protocol is tuned to quickly increase its transmission rate without the need for a long SS. Once arriving at the output (B), another TCP connection is used to transmit the packets to the destination. In a satellite environment, the long delay link may lead directly to the destination, so another TCP connection is not required. Because packets have already been acknowledged, any loss between the input of the link (A) and the destination must be locally retransmitted on behalf the source. Also, ACKs from the receiver must be discarded silently (at B) to not confuse the source. This approach is called TCP spoofing.

The main gain in performance comes from not using SS on the long delay link. The window increases quickly, which improves performance; but spoofing still has many drawbacks. First, it breaks the end-to-end semantics of TCP; a packet is acknowledged before reaching its destination. Also, it doesn't work when encryption is accomplished at the IP layer, and it introduces a heavy overload on network routers. Further, the transfer is vulnerable to path changes, and symmetric paths are required to be able to discard the ACKs before they reach the source. Spoofing can be seen as a particular solution to some long delay links. It is interesting when the long delay link is the last hop to the destination. This solution is often used in networks providing high-speed access to the Internet via geostationary earth orbit (GEO) satellite links.

#### SOLUTIONS TO UNFAIRNESS

Two trends exist to improve the fairness of TCP. The first tries to solve the problem at the TCP level by accelerating the window growth for long RTT connections. The second tries to solve the problem inside the network. Some intelligence is proposed to be added in network nodes in order to allocate the bandwidth fairly among different connections.

As an example of a TCP-level solution, we use the proposition in [6], the constant rate algorithm. The window is increased in CA by a factor inversely proportional to  $RTT^2$ . The result is a constant increase rate of the throughput regardless of RTT, thus better fairness. The first problem in this proposition is the choice of the increase rate. Also, accelerating window growth while preserving the ACK clock results in large bursts for long RTT connections.

Inside the network, fairness is improved by isolating the different connections from each other. Given that congestion control in TCP is based on losses, isolation means that a congested node must manage its buffer intelligently to distribute drops on the different connections in such a way that they get the same throughput. Many buffer management policies have been proposed. Some of these policies, such as RED [10], drop incoming packets with a certain probability when the queue length or its average exceed a certain threshold. This distributes losses on the different connections proportionally to their throughput without requiring any per-connection state. However, it has been shown in

[14] that dropping packets in proportion to the throughput doesn't always lead to fairness, especially if the bottleneck is crossed by unresponsive traffic. With a first-in first-out (FIFO) scheduler, the connection share of the bandwidth is proportional to its share of the buffer. Better fairness requires control of the buffer occupancy of each connection. We find here another set of policies, as Flow RED [14], which try to improve fairness by sharing the buffer space fairly between active connections. This ensures that each connection has at least a certain number of places in the queue, which isolates connections sending at small rates from aggressive ones. This improves fairness, but at the same time increases buffer management overhead over a general drop policy such as RED.

Solving the fairness problem at the TCP level has the advantage of keeping routers simple, but it is not enough given the prevalence of non-TCP-friendly traffic. Some mechanisms in network nodes are required to protect conservative TCP flows from aggressive ones. Network mechanisms are also required to ensure fairness at a level below or above TCP, say at the user or application level. A user (e.g., running XFTP) may cheat and establish many TCP connections in order to increase its share of the bandwidth. The packets generated by this user must be considered by the network as a single flow. This requires an aggregation in flows of TCP connections. The level of aggregation determines the level of fairness we want.

#### NONCONGESTION LOSSES

TCP considers the loss of packets as a signal of network congestion and reduces its window consequently. This results in severe throughput deterioration when packets are lost for reasons other than congestion. Noncongestion losses are mostly caused by transmission errors. A packet may be corrupted while crossing a poor-quality (e.g., wireless) link. The solutions proposed to this problem can be divided into two main categories. The first consists in hiding the lossy parts of the Internet so that only congestion losses are detected at the source. The second type of solution consists in enhancing TCP with some mechanisms to help it to distinguish between different types of losses.

#### HIDING NONCONGESTION LOSSES

Noncongestion losses are recovered locally without the intervention of the source. This can be accomplished at the link or TCP level. A link is any transmission medium located between two adjacent routers. The objective is to approach the typical network to which TCP is tuned: one formed by the interconnection of good-quality transmission media.

Link-Level Solutions — Two well-known mechanisms exist for the improvement of link quality: automatic repeat request (ARQ) and forward error correction (FEC). ARQ is efficient when losses are not frequent and propagation delay is not important. Extra bandwidth is consumed only when a packet is retransmitted. However, ARQ may interfere with TCP mecha-

Two trends exist to improve the fairness of TCP.

The first tries to solve the problem at the TCP level by accelerating the window growth for long RTT connections. The second tries to solve the problem inside the network.

对于noncongestion losses, 处理的方法分为两类:

1. 隐藏, 让tcp层感觉不到丢包
2. 增强, 让tcp分辨出拥塞丢包和 poor-quality丢包



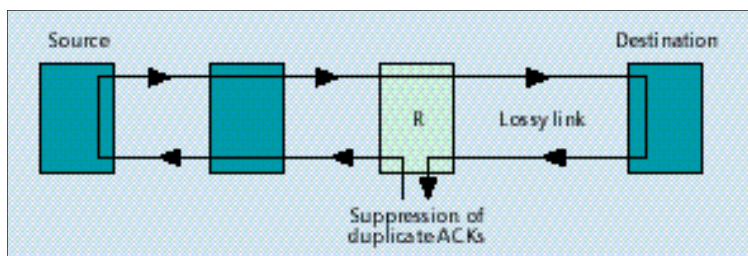


Figure 3. A TCP-aware link layer.

nisms [2]. If the link layer doesn't provide in-sequence delivery of packets, TCP packets following the loss keep arriving at the destination and trigger the transmission of duplicate ACKs. These duplicate ACKs reach the source while the link layer is retransmitting the packet. This causes an unnecessary window reduction which we try to avoid. **The proposed solution to this problem was to use a TCP-aware ARQ protocol (Fig. 3).** The link layer suppresses the duplicate ACKs (at R) so that they don't reach the source. If the link layer fails to retransmit the packet, the source will timeout and retransmit the packet itself. **This solution is applicable only when the lossy link is the last hop to the destination.** If the lossy link is followed by other routers, congestion losses will be hidden, which must be avoided.

FEC consists of sending some redundant information in order to rebuild the corrupted part of the packet. The drawback of this technique is that the redundant information is not used when the link is in the good state. Thus, it presents a certain waste of bandwidth. Also, the computation of this redundant information requires extra CPU processing time, memory, and blocking delay. However, the advantages of FEC are worth the cost. First, corrupted packets are corrected directly without retransmission, **which is important for long delay lossy links such as satellites ones.** Also, FEC doesn't interfere with TCP mechanisms. For these reasons, this technique has been recommended in a satellite environment [13]. Convolutional coding, Viterbi decoding together with interleaving techniques, and Reed-Solomon encoding are widely used to render satellite links as clean as terrestrial ones.

**TCP-Level Solutions** — These solutions try to improve link quality by retransmitting packets at the TCP level rather than at the link level. **A TCP agent** in the router at the input of the lossy link keeps a copy of every data packet. It discards this copy when it sees the ACK of the packet, and **it retransmits the packet on behalf of the source when it detects a loss.** This technique has been proposed for terrestrial wireless networks where the delay is not so important as to require the use of FEC. The TCP agent is placed in the base station at the entry of the wireless network. Two possible implementations of this agent exist.

The first implementation (e.g., indirect TCP [2]) consists of terminating the originating TCP connection at the entry of the lossy link. The agent acknowledges the packets and takes care

of handing them to the destination. A TCP connection well tuned to a lossy environment (e.g., TCP-SACK [9]) can be established across the lossy network. A different transport protocol can also be used. This solution breaks the end-to-end semantics of the Internet. Also, it causes difficulties during handover since a large state must be transferred between base stations.

The second implementation (Snoop protocol [2]) respects the end-to-end semantics. The intermediate agent doesn't terminate the TCP connection; it just keeps copies of data packets and doesn't generate any artificial ACK. Nonduplicate ACKs sent by the destination are forwarded to the source. Duplicate ACKs are stopped. A packet is retransmitted locally when three duplicate ACKs are received or a local timeout expires. This local timeout is set, of course, to a value less than that of the source. As in the link-level case, interference may happen between the source and agent mechanisms. In fact, this solution is no other than link-level recovery implemented at the TCP level. **Again, because it hides all losses, congestion losses must not occur between the Snoop agent and the destination.**

## END-TO-END SOLUTIONS

Improvement of the quality of Internet links will not hide all noncongestion losses. First, even if they are concentrated on some links, such losses may appear anywhere in the network, and it seems impossible to track them all. Second, local mechanisms such as FEC may fail to recover from lost packets. The addition of some end-to-end mechanisms to improve TCP reaction to noncongestion losses should further improve performance.

Two approaches exist in the literature. ~~The first consists of explicitly informing the source of the occurrence of a noncongestion loss via an Explicit Loss Notification (ELN) signal [6].~~ **The source reacts by retransmitting the lost packet without reducing its window.** An identical signal has been proposed to halt congestion control at the source when a disconnection appears due to handover in a cellular network. The difficulty with such a solution is that a packet corrupted at the link level is discarded before reaching TCP, and then it is difficult to get this information.

The second approach is to improve the congestion control provided by TCP rather than recovery from noncongestion losses. We mention it here because it consists a step toward a solution to the problem of losses on an end-to-end basis. The proposed solutions aim to **decouple congestion detection from losses.** With some additional mechanisms in the network or at the source, the congestion is detected and the throughput reduced before the overflow of network buffers. **As solutions, we find the Vegas version of TCP [15] and the Explicit Congestion Notification (ECN) proposal [12].** In Vegas, the RTT of the connection and the window size are used to compute the number of packets in network buffers. The window is decreased when this number exceeds a certain threshold. With ECN, an explicit signal is sent by the routers to indicate congestion to TCP sources rather than dropping packets.

If all the sources, receivers, and routers are

compliant (according to Vegas or ECN), congestion losses will considerably decrease. The remaining losses could be considered to be caused mostly by problems other than congestion. Given that noncongestion losses require only retransmission without window reduction, the disappearance of congestion losses may lead to the definition at the source of a new congestion control algorithm which reacts less severely to losses.

This ideal behavior doesn't exist in today's networks. In the absence of any feedback from the network as with Vegas, the congestion detection mechanism at the source may fail; here, congestion losses are unavoidable. If the source bases its congestion control on explicit information from the network as with ECN, some non-compliant routers will not provide the source with the required information, dropping packets instead. A reduction of the window is necessary in this case. For these reasons, these solutions still consider losses as congestion signals and reduce their windows consequently.

## BANDWIDTH ASYMMETRY

TCP uses the ACK clock to predict what is happening inside the network. It assumes implicitly that the reverse channel has enough bandwidth that ACKs traverse it without being disturbed. This was almost true with the so-called symmetric networks where the forward and the reverse directions have the same bandwidth. However, some of today's networks (e.g., direct broadcast satellite and asymmetric digital subscriber loop networks) tend to increase capacity in the forward direction using satellite links or cables, whereas a low-speed channel such as a dialup modem line is used to carry ACKs back to the source (Fig. 4). Even if ACKs are smaller in size than data packets, the reverse channel is unable to carry the high rate of ACKs. The result is congestion and losses on the ACK channel (at A).

The congestion in A increases the RTT of the connection and causes loss of ACKs. The increase in RTT reduces throughput and increases end-to-end delay. Also, it slows window growth, which further impairs performance when operating on a long delay path or in a lossy environment.

The loss of ACKs disturbs one of the main functionalities of the ACK clock: smoothing the transmission. The window slides quickly upon receipt of an ACK covering multiple lost ACKs, and a burst of packets is sent, which may overwhelm the network buffers in the forward direction [1, 8]. Also, the loss of ACKs slows down the growth of the congestion window, which results in poor performance for long delay paths and lossy links.

The proposed solutions to this problem can be divided into receiver-side solutions, which try to solve the problem by reducing the congestion on the return path, and source-side solutions, which try to reduce TCP burstiness.

The first receiver-side solution is to compress the headers of TCP/IP packets on a slow channel (A-B) to increase its capacity in terms of ACKs per unit of time (e.g., SLIP header compression [12]). It profits from the fact that most of the information in a TCP/IP header doesn't change during the connection lifetime.

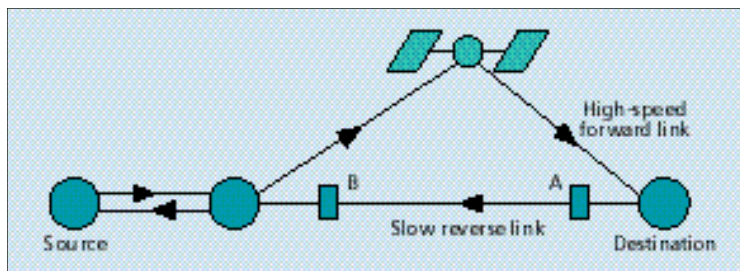


Figure 4. An asymmetric path.

The other solutions propose to reduce the rate of ACKs to avoid congestion. The first proposition is to delay ACKs at the destination [1]. An ACK is sent every  $d$  packets, and an adaptive mechanism has been proposed to change  $d$  as a function of the congestion on the ACK path. Another proposition [1] keeps the destination unchanged and filters ACKs at A. When an ACK arrives, the buffer is scanned to see if another ACK (or a certain number of ACKs) of the same connection is buffered. If so, the new ACK is substituted for the old one. ACKs are filtered to match their rates to the rate of the reverse channel. Normally, in the absence of artificial filtering, ACKs are filtered sometime later when the buffer gets full. The advantage of this solution is that the filtering is accomplished before the increase in RTT.

Solutions at the sender side are proposed to reduce the burstiness of TCP. Note that this problem is caused by the reliance of TCP on the ACK clock, and we believe it cannot be completely solved without any kind of packet spacing. First, a limit on the size of bursts sent by TCP has been proposed. However, with systematic loss of ACKs, limiting the size of bursts limits the throughput of the connection. Second, it has been proposed in [1] to reconstruct the ACK clock at the output of the slow channel (at B). When an ACK arrives at B, all the missing ACKs are generated, spaced by a time interval derived from the average rate at which ACKs leave the slow channel. This reconstruction may contain a solution to this particular problem. However, the general problem of TCP burstiness upon loss of ACKs still exists.

## CONCLUSIONS

In this work, we summarize most of the problems that face TCP in the Internet of today. We present the impact of the new characteristics of Internet paths on TCP performance, mainly the bandwidth-delay product, round-trip time, non-congestion losses, and bandwidth asymmetry. The study is accomplished independent of network type. We outline some of the solutions proposed in the literature. We show the utility of these solutions on a general path having a certain characteristic. We believe that the two main problems of TCP to be solved are burstiness and the coupling between congestion detection and error control. Solutions to these problems are required to make TCP able to cope with the heterogeneity of today's transmission media on an end-to-end basis.

We believe the two main problems of TCP to be solved are burstiness and the coupling between congestion detection and error control. Solutions to these problems are required to make TCP able to cope with the heterogeneity of today's transmission media on an end-to-end basis.

## REFERENCES

- [1] H. Balakrishnan, V. Padmanabhan, and R. Katz, "The Effects of Asymmetry on TCP Performance," ACM MOBICOM, Sept. 1997.
- [2] H. Balakrishnan et al., "A Comparison of Mechanisms for Improving TCP Performance over Wireless Links," ACM SIGCOMM, Aug. 1996.
- [3] T. Henderson and R.H. Katz, "Transport Protocols for Internet-Compatible Satellite Networks," IEEE JSAC, Feb. 1999.
- [4] V. Jacobson, "Congestion Avoidance and Control," ACM SIGCOMM, Aug. 1988.
- [5] C. Barakat and E. Altman, "Analysis of TCP with Several Bottleneck Nodes," IEEE GLOBECOM, Dec. 1999.
- [6] S. Floyd, "Connections with Multiple Congested Gateways in Packet-Switched Networks Part 1: One-Way Traffic," Comp. Commun. Rev., Oct. 1991.
- [7] T. V. Lakshman and U. Madhow, "The Performance of TCP/IP for Networks with High Bandwidth-Delay Products and Random Loss," IEEE/ACM Trans. Networking, June 1997.
- [8] T. V. Lakshman, U. Madhow, and B. Suter, "Window-Based Error Recovery and Flow Control with a Slow Acknowledgment Channel: A Study of TCP/IP Performance," IEEE INFOCOM, Apr. 1997.
- [9] K. Fall and S. Floyd, "Simulation-Based Comparisons of Tahoe, Reno, and SACK TCP," Comp. Commun. Rev., July 1996.
- [10] S. Floyd and V. Jacobson, "Random Early Detection Gateways for Congestion Avoidance," IEEE/ACM Trans. Networking, Aug. 1993.
- [11] J. Hoe, "Improving the Start-Up Behavior of a Congestion Control Scheme for TCP," ACM SIGCOMM, Aug. 1996.
- [12] M. Allman et al., "Ongoing TCP Research Related to Satellites," Internet draft, work in progress, Sept. 1999.
- [13] M. Allman, D. Glover, and L. Sanchez, "Enhancing TCP over Satellite Channels Using Standard Mechanisms," RFC 2488, Jan. 1999.
- [14] D. Lin and R. Morris, "Dynamics of Random Early Detection," ACM SIGCOMM, Sept. 1997.
- [15] L. Brakmo and L. Peterson, "TCP Vegas: End to End Congestion Avoidance on a Global Internet," IEEE JSAC, Oct. 1995.

## BIOGRAPHIES

CHADI BARAKAT (cbarakat@sophia.inria.fr) received his diploma degree in electrical and electronics engineering from the Engineering Faculty of Lebanese University, Beirut, in 1997. In 1998 he obtained the D.E.A. degree in networks and distributed systems from the University of Nice, Sophia Antipolis, France. Since 1998 he has been working toward a Ph.D. degree within the MISTRAL team at National Research Institute for Computer Science and Control (INRIA), Sophia Antipolis. His main research topics are congestion and error control in computer networks, performance evaluation of communication protocols, and integration of new transmission media such as satellite networks into the Internet.

EITAN ALTMAN (altman@sophia.inria.fr) received a B.Sc. degree in electrical engineering (1984), a B.A. degree in physics (1984), and a Ph.D. degree in electrical engineering (1990), all from the Technion-Israel Institute, Haifa. In 1990 he further received his B.Mus. degree in music composition from Tel Aviv University. Since 1990 he has been with INRIA, Sophia Antipolis, France. His current research interests include performance evaluation and control of telecommunication networks, stochastic control, and dynamic games. In recent years he has applied control theoretical techniques in several joint projects with the French telecommunications company, France Telecom.

WALID DABBOUS (dabbous@sophia.inria.fr) graduated from the Faculty of Engineering of the Lebanese University, Beirut, in 1986. He obtained his D.E.A. and Doctorat d'Université from the University of Paris XI in 1987 and 1991, respectively. He joined the RODEO Team within INRIA in 1987. He has been a staff researcher at INRIA since 1991, and a leader of the RODEO team since 1996. His main research topics are high-performance communication protocols, congestion control, reliable multicast protocols, audio and video conferencing over the Internet, efficient and flexible protocol architecture design, and the integration of new transmission media such as satellite links in the Internet. He is co-chair of the UDLR working group of the Internet Engineering Task Force.