

PSockets: The Case for Application-level Network Striping for Data Intensive Applications using High Speed Wide Area Networks

H. Sivakumar
sharinat@eecs.uic.edu
University of Illinois at Chicago
Chicago, IL, USA

S. Bailey
sbailey@infoblox.com
Infoblox Inc.
Evanston, IL, USA

R. L. Grossman
grossman@uic.edu
University of Illinois at Chicago
Chicago, IL, USA
Magnify Inc.
Chicago, IL, USA

Abstract

Transmission Control Protocol (TCP) is used by various applications to achieve reliable data transfer. TCP was originally designed for unreliable networks. With the emergence of high-speed wide area networks various improvements have been applied to TCP to reduce latency and achieve improved bandwidth. The improvement is achieved by having system administrators tune the network and can take a considerable amount of time. This paper introduces PSockets (Parallel Sockets), a library that achieves an equivalent performance without manual tuning. The basic idea behind PSockets is to exploit network striping. By network striping we mean striping partitioned data across several open sockets. We describe experimental studies using PSockets over the Abilene network. We show in particular that network striping using PSockets is effective for high performance data intensive computing applications using geographically distributed data.

1. INTRODUCTION

With the rapid advancements in networking technologies, high-speed networks such as Abilene and vBNS (very high speed backbone network service) are becoming more common. Although distributed data mining and data intensive computing applications have an urgent need for the throughput provided by these high performance networks, achieving the desired throughput in practice can sometimes be difficult due to the network tuning that is usually required to achieve optimal performance. The key focus of this paper is to help applications to achieve the best end-to-end performance (maximum throughput) over high-speed wide area networks without modifications to the existing network.

Most data applications need reliable data transfer and hence utilize the Transmission Control Protocol (TCP) [1]. TCP was originally designed for unreliable networks that required reliable data transfers between the source and destination machines. A common way to improve the bandwidth and reduce the latency while using TCP on reliable high-speed wide area networks is to tune the TCP window size appropriately [2, 7]. This involves the system administrator at both ends and it can be time consuming to find the best window size for optimal performance.

In this paper, we introduce a C++ library called PSockets (Parallel Sockets) which applications can use without the

necessity of tuning the TCP window size and yet still achieve near optimal utilization of the network bandwidth. The ¹idea is a simple one and an old one – we divide the data into partitions, open several sockets, and stripe the data over the sockets. We call this approach *network striping*. This can be broadly viewed as analogous to striping data over arrays of disks [9].

We believe that this paper makes the following contributions:

1. We describe five experiments demonstrating that network striping is an effective mechanism for applications to achieve high throughput on wide area high performance networks.
2. We describe a C++ library we have developed called PSockets that allows applications to quickly incorporate network striping with a few simple function calls.
3. We demonstrate that data intensive computing applications can easily incorporate PSockets into their design and transparently exploit wide area high performance networks.

We believe that this paper is novel for the following reasons:

1. With the exception of [8], of which we have just recently become aware, most approaches to improving performance have focused on system level tuning of TCP window size, not application-level network striping. In this paper, as in [8], we show that this alternative approach is equally effective, and in certain situations, more effective, than tuning window size.
2. Prior work on network tuning has focused largely on network measurements and not on application level measurements. In this paper, we also show that our approach using network striping allows the development of effective data intensive computing applications, an area of rapidly increasing importance.
3. We provide experimental data which shows that in certain situations network striping using PSockets is more effective than tuning window size, even when theoretically this should not be expected. For example, it is common that as the complexity of a network grows it is more likely for there to be some hardware or software

¹ 0-7803-9802-5/2000/\$10.00 (c) 2000 IEEE.

configuration errors somewhere along the path connecting the sender and receiver, which will limit network performance. This occurred quite often in our experiments. In this case, we have observed that network striping at the application level is still an effective mechanism for increasing network throughput, while; on the other hand, tuning window size is not nearly as effective in practice.

This paper is organized as follows. Pockets is described in Section 3. Section 4 discusses four sets of experiments conducted over the Abilene network. Section 5 discusses some of the applications we have built using Pockets.

2. BACKGROUND

TCP was designed to operate reliably over almost any transmission medium regardless of transmission rate, delay, corruption, duplication, or reordering of segments [1]. TCP implementations currently adapt to transfer rates in the range of 100 bps to 10^7 bps and round-trip delay in the range of 1ms to 100 seconds [2]. TCP guarantees reliable in-order delivery of data sent from the source to the destination. Since many applications need reliable data delivery, TCP has been very widely used.

Typically an application requests data to be sent to another application at a remote machine. The TCP stack in the kernel of an operating system handles requests from various applications and passes data to the network. TCP partitions the data into *segments* to be sent over appropriate transmission media and waits for an acknowledgement from the receiver to know that a particular segment has been received correctly. TCP achieves this by having windows on both the sender and the receiver that can throttle the rate at which data can be sent [10].

Due to the recent advancements in fiber optics and related technologies, networks with very high transmission speeds are becoming more common. TCP was not designed for these types of networks [2] and recently there has been research addressing this issue [2, 3, 4, 5, 11]. TCP window size is a key tuning factor for networks with large delays between the sender and the receiver. RFC1323 [2] provides details on achieving better TCP performance over such networks. Automatic TCP Buffer tuning to achieve increased bandwidth on a socket level is detailed in [3]. Improvements on the TCP algorithm to transmit data quickly and efficiently are explained in [4] and [5]. TCP performance improvements over satellite links were carried out in [8] which involved large delays. TCP performance depends not upon the transfer rate itself, but rather upon the product of the transfer rate and the round-trip delay. This “bandwidth*delay product” measures the amount of data that would fill the pipe. It is the buffer space required at the sender and the receiver to obtain the maximum throughput over the established TCP connection, i.e., the amount of unacknowledged data that TCP must handle in order to keep the pipeline full. TCP

performance problems arise when the bandwidth*delay product is large. This is referred to as “long, fat pipe” and a network containing this path is a long fat network (LFN). The three fundamental performance problems with the current TCP over LFN paths are the window size limit, the recovery from losses and the round-trip measurements. Several RFC have been developed and implemented to overcome these issues and achieve maximum performance with the current TCP over LFN.

The key factor in achieving maximum performance from TCP over LFN is the TCP window size. The TCP window size has to be tuned at both the source and the destination in order to achieve the maximum throughput. This involves changing parameters for the TCP stack in the kernel and has to be done by the system administrators at both ends [6]. Also, after the tuning, performance analysis has to be done to fine tune this window size. Typically this operation can take anywhere from a few hours to even weeks [6]. In this paper we propose an alternative method involving striping data over several sockets. Applications can use the Pockets library and don't have to bother tuning the window size. Our library achieves equivalent performance to the one that would be achieved while setting the best window size.

A similar method with satellite applications in mind was proposed in [8]. Both use several open sockets and the basic principle is the same. In [8], the data is divided into 8KB blocks and the blocks are striped across several sockets. In Pockets, the segmentation and re-assembly of the data across various sockets is done with the help of pointers in C++. Striping across open sockets as 8KB blocks involves the overhead of keeping track of data on different sockets and then re-assembling them at the destination. We had minimal overhead by striping the data equally across various sockets as compared to striping as 8KB blocks. Hence, we have found it to be more advantageous to divide the data into equal partitions and send the partitions over the open sockets asynchronously. In this way we incur very little overhead in re-assembly and achieve better performance.

3. SOCKETS

When an application needs to transfer data reliably over a network, it opens up a TCP/IP socket connection between the sender and the receiver. When the sender and the receiver are connected by a high speed wide area network the bandwidth*delay product is quite high (greater than 64KB). In order to achieve the maximum TCP Performance on such networks the sender and the receiver's TCP window size is set to the bandwidth*delay product [2]. Enabling RFC 1323 on a system has to be done by the system administrators at both ends and not by the application. This could typically take a few days or sometimes even weeks. Therefore, in practice most applications are limited to the default maximum window size.

Using Pockets, individual applications can overcome this limitation by using multiple sockets. Since the limitation of the TCP window size is only on a single socket, Pockets opens multiple socket connections between the sender and the receiver. Pockets then divides the data to be transferred into equal partitions, where the number of open sockets determines the number of partitions. Pockets then sends the data partitions asynchronously over the multiple socket connections. The TCP/IP stack in the kernel multiplexes the data on the various connections so that, in effect, the partitions appear to be sent in parallel on the multiple connections. In this way, Pockets is able to achieve the maximum data transfer rate possible on a high-speed wide area network without any network tuning. We call this approach *network striping*.

Figure 1 shows TCP packets between a source and destination pair at a given instant of time over a network. Assuming that the default TCP window size is set to 64KB and the bandwidth*delay product is equal to 192KB, Figure 1 shows three scenarios for a LFN. In each case 192KB of data is to be transmitted from source to destination.

- Figure 1 (a): Here, the connection between the source and the destination has the default window size (64KB). At the source, the first 64KB of data is sent since the TCP window size is 64KB. The sender then waits for an acknowledgement before sending any more data. Since the source waits for an acknowledgement from the destination the connection pipe between the source and destination is not completely filled. Therefore the maximum throughput attainable over the connection is not achieved. This limitation is overcome by using RFC 1323.
- Figure 1 (b): This diagram shows the connection between the source and the destination after enabling RFC 1323, which recommends the use of large TCP window sizes for connections having large bandwidth*delay products. In our example the TCP window size is set to 192KB (bandwidth*delay product). In this scenario, the source sends the entire 192KB of data to be transmitted to the destination without waiting for an acknowledgement. This achieves the maximum data rate between the source and the destination. The entire pipe is filled with segments up to the TCP window size.
- Figure 1 (c) shows our method using multiple sockets. No network tuning is required in this scenario. The TCP window size is left at the default value of 64KB. In this example, three sockets are used to fill up the pipe since each socket has a window size of 64KB. This achieves the same performance as in Figure 1 (b).

Pockets uses application level network tuning to achieve the maximum throughput between a source and a destination. The throughput attained by Pockets is equivalent to the one obtained over a finely tuned network which has a large bandwidth*delay product. If the TCP driver does not use a

good recovery algorithm, then quite a few packets may be dropped and have to be retransmitted. This reduces the data rate of the connection. When the TCP window size is set to the bandwidth*delay product for the high-speed wide area network, the retransmission of the transition data (the data which has been sent from the source and not yet arrived at the destination) can be costly and will affect the performance. By using Pockets, we can use a smaller TCP window size, reducing the number of packets which need to be retransmitted

Using the Pockets library, applications don't have to tune for the best TCP window size to achieve the maximum throughput. The Pockets library hides information on the number of open sockets, the segmentation, and the re-assembly of the data. Pockets uses an API similar to the standard socket "send" and "receive".

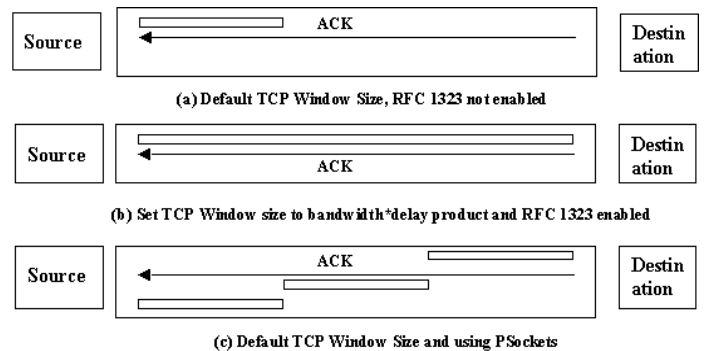


Figure 1: TCP packets at a given time over a LFN

4. EXPERIMENTAL RESULTS

	Machine at Chicago, IL	Machine at Ann Arbor, MI
Processor	Dual processor Pentium II 450 MHz	Uni processor Pentium II 450 MHz
Network card	100 Mb/s fast Ethernet card	100 Mb/s fast Ethernet card
Operating System	Red Hat Linux 6.1	Red Hat Linux 6.1
External connectivity	ATM Switch	Fast Ethernet Switch

Table 1: Configuration of the Machines used for experiments

Experiments were conducted between machines at Chicago, IL and Ann Arbor, MI. Table 1 describes the configuration of the machines. The connectivity between Chicago and Ann Arbor is an OC3 line (155 Mb/s). This link was limited to an effective rate of 100 Mb/s due to a fast Ethernet switch in the path with this limitation.

We performed several experiments described below to study the performance of Pockets. In all the experiments

performed, the overhead to segment and re-assemble the data has been included.

a) Varying number of sockets: The first test conducted was to vary the number of parallel sockets used for data transfer. The goal of the experiment was to find the optimal number of sockets for a particular buffer size. Figure 2 shows the results of our experiment for the TCP default window size of 64KB. The amount of data transmitted in each experiment was 64KB. It was observed that the maximum throughput was obtained for a PSocket size of 12. Other experiments conducted (with a constant TCP window size and constant amount of data transfer) revealed that the maximum throughput was attained using PSocket sizes between 6 and 16. Using PSockets we were able to achieve a maximum throughput of 76 Mb/s of application data (not including other lower level overheads).

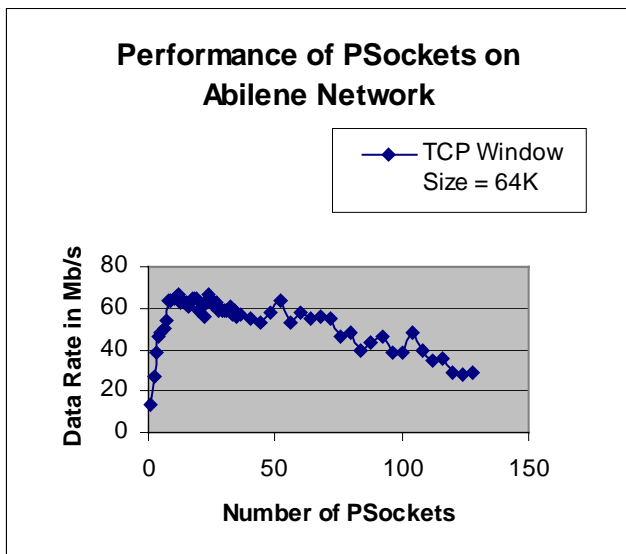


Figure 2: Performance of PSockets

b) Varying the amount of data transfer: In this experiment we varied the amount of data transferred. We kept the TCP window size set to the default maximum of 64KB. To compare the performance of PSockets against standard software, we utilized Iperf [6], a network performance tool developed by NLANR. Figure 3 shows the results of our experiments for data transfer amounts varying from 1KB to 4MB using PSockets as well as Iperf. Observe that the throughput increases with an increase in the number of sockets for a particular amount of data and then finally drops down. See Figure 1. For a particular PSocket size the increase in the size of data transferred from 1KB to 4MB did not have a noticeable effect on the throughput obtained. We concluded that for a fixed number of sockets, varying the amount of data transferred did not affect the throughput.

c) Performance of PSockets with and without network tuning: The third test we performed was to vary the TCP window size and see the performance of using RFC 1323 as

well as PSockets. Figure 4 shows the results for two TCP window sizes, 64KB (default TCP window size on Linux operating system) and 256KB (A TCP window size closer to the bandwidth*delay product between the machines used for the experiment). The data transmitted in each experiment was 64KB. This clearly validates our argument that applications using PSockets can obtain throughputs equivalent to those obtained after careful network tuning (RFC 1323) of the high-speed wide area networks. Observe that we are able to achieve a throughput of 57.7 Mb/s without network tuning and 67.7 with network tuning. Note that the maximum throughput attained using Iperf for a TCP window size of 256KB was only 47.7 Mb/s. This could potentially be a problem in the implementation of the TCP/IP driver or the network itself [7].

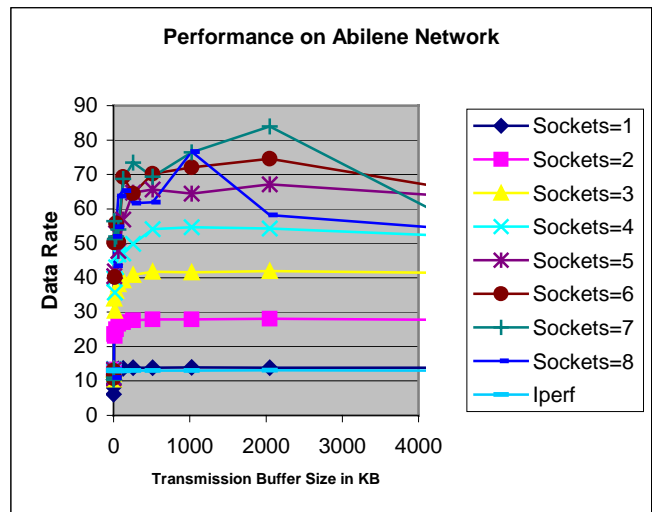


Figure 3: Performance results using PSockets on Abilene network for various transmission Buffer sizes

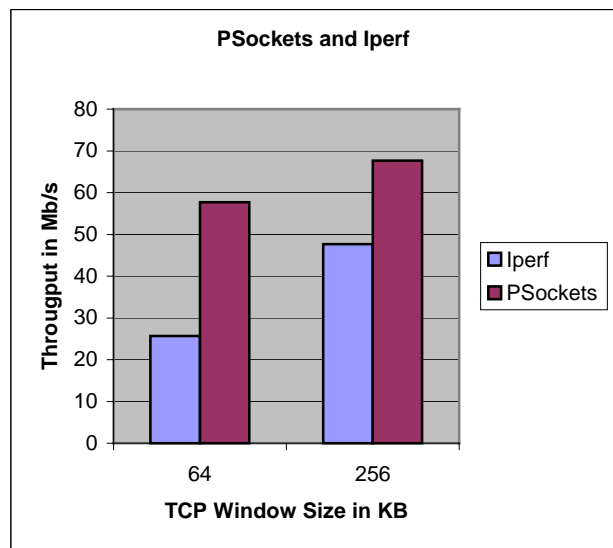


Figure 4: PSockets without and with network tuning

d) Performance of PSockets with various TCP window sizes: Figure 5 shows the performance of PSockets (8 sockets) for TCP window sizes 1KB to 4MB. The bandwidth*delay product for the connection is calculated to be 245KB. From our second experiment we concluded that the amount of data transferred did not affect the throughput obtained using PSockets. Therefore we used 8KB as the data transfer unit for this experiment. It can be noticed from the graph that the maximum throughput was obtained for the TCP window size 512KB. Therefore, as mentioned in [7], setting the TCP window size equal to the bandwidth*delay product actually does not give the best performance. The PSockets performance is compared against Iperf for the same window size. Observe that the difference in throughput obtained using PSockets and Iperf narrowed as the TCP window size increases. Since PSockets uses 8 sockets to send data, it performs better than Iperf up to the TCP window size of 8KB. Beyond the 8KB TCP window size, one would expect both PSockets and Iperf to be able to achieve the same throughput. Since the Iperf performance was always less than that of PSockets we tried to verify our suspicion that there was a potential problem with the TCP/IP driver or the network itself. Iperf will not achieve the maximum throughput with a single client when there are problems with either the TCP/IP driver or the network [7]. When Iperf was simulated with two clients connecting to the server at the same time (similar to running PSockets with a value 2), this defect was confirmed. Iperf with 2 client connections at the same time was able to achieve an aggregate throughput more than that of Iperf with a single client.

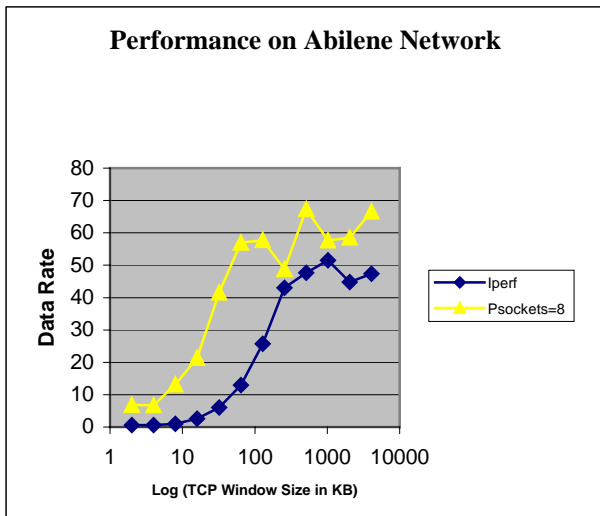


Figure 5: Performance Comparison of PSockets and Iperf with varying TCP Window Sizes

e) Varying the number of clients. The next experiment we performed was to increase the number of applications running PSockets. We used a PSocket size of 8 for this experiment and a TCP window size of 64KB and the amount of data transferred by each application was 2MB. Observe that the aggregate throughput obtained by applications using

PSockets increased as the number of applications increased from 1 to 7 and then gradually dropped down for 8 or more. This is due to the fact that the driver has to handle 64 sockets. This can be observed from Figure 6. It was observed that the throughput obtained by each application was nearly the same indicating that each application using PSockets was given an equal share of the bandwidth.

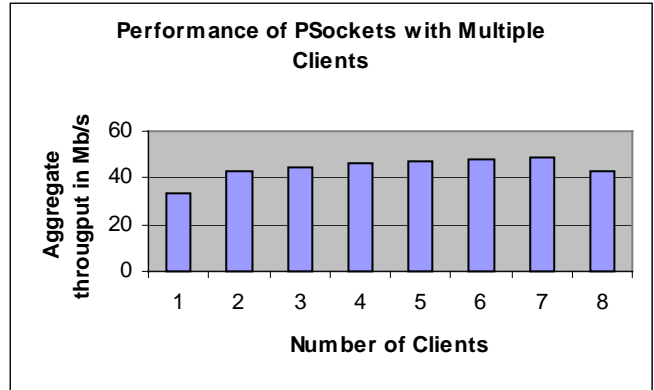


Figure 6: Performance of PSockets with Multiple Clients

Additional experiments were conducted to compare the results of applications not using PSockets with applications using PSockets. An application using PSockets (8 sockets) was run along with an application using a single socket. The maximum attainable throughput was equally divided among the 9 sockets. We can conclude that the applications not using PSockets will get lower throughput when run along with applications using PSockets. We analyzed the packet drops in three applications, which transferred 512KB of data – one using PSockets (8 sockets) with the default window size, the second application using one socket with the default TCP window size, and the third application using one socket, but with the TCP window size set to the bandwidth*delay product. The packet drop percentages in the three applications were 0.03%, 0.02% and 0.07% respectively. The software tools *tcpdump* and *tcptrace* [11] were used to compute this. The bandwidth attained in the three applications was 57 Mb/s, 38 Mb/s and 59 Mb/s respectively. This indicates that there is more contention for the available bandwidth while using PSockets with the default TCP window size than while using a single socket with the same TCP window size. On the other hand, we are able to achieve a very high throughput while using PSockets without the cost of network tuning. Note that the bandwidth achieved by an application with a single socket along with the TCP window set to bandwidth*delay product is slightly greater than that achieved by PSockets but the packet drop percentage is nearly twice as much as that of the application using PSockets. This test clearly confirms that applications using PSockets without network tuning will be able to achieve equivalent bandwidth to that of a well tuned network.

The slow start and the congestion avoidance algorithms [5] for a wide area network minimize the packet loss during congestion. All our experiments were run on a Linux system,

and the TCP/IP driver did not have options to configure the slow start algorithm. Hence the effect of slow start algorithm while using Pockets has not been studied in this paper.

From all of the above-mentioned experiments, we conclude that Pockets helps in achieving a throughput closer to the maximum throughput attainable over the network by using application level tuning rather than the more time-consuming network tuning. Pockets would be able to help in achieving better throughput even when there are inherent problems with the driver or the network. Applications using Pockets will be able to extract a greater bandwidth share over a high-speed network while competing against other normal applications.

5. APPLICATIONS

We built a geographically distributed data intensive computing application with Pockets to mine high energy physics data. The data was located at Chicago, IL, Ann Arbor, MI and Arlington, VA. The Arlington link was limited due to a DS3 (45Mb/s), while the Ann Arbor was limited due to a fast ethernet switch of 100Mb/s. The Chicago machines were limited due to the OC3 link (155 Mb/s). Six Linux machines pulled data from these locations into Portland, OR at the SuperComputing '99 conference floor. We were able to obtain a maximum throughput of around 30-35Mb/s while using the Physics application since the data was striped across 6 machines (2 machines at each site).

In addition, we developed a simple application benchmark program to test the raw performance of Pockets. We were able to achieve a maximum throughput of around 240Mb/s using this benchmark. The maximum attainable throughput was equal to around 300 Mb/s. The input link at the floor was an OC12 (622 Mb/s).

Currently we are using the Pockets library to build a high performance wide area data storage application called Osiris. Osiris stripes row-column data across various nodes distributed geographically and interconnected by a high speed wide area network. Osiris is aimed at applications that need large volumes of data to be retrieved quickly. Pockets is currently available for the Linux platform and can be downloaded from <http://www.ncdm.uic.edu>.

6. CONCLUSIONS

We have developed a library called Pockets, which helps wide area applications that need to move large amounts of

data. Even though Pockets achieves a performance equivalent to that obtained with RFC 1323 enabled, it is much easier to use since no tuning is required. Typically, tuning the network can be quite labor intensive, since it requires work by system administrators on both ends. With the Pockets library, developers need not worry about this. Since Pockets has the same API as that of regular sockets it is very easy for application developers to use.

REFERENCES

- [1] J. Postel, "Transmission Control Protocol", Request for Comments 0793, Sep. 1981.
- [2] Van Jacobson, Robert Braden, and Dave Borman. "TCP extensions for high performance", Request for Comments 1323, May 1992.
- [3] Jeffery Semke, Jamshid Mahdavi and Matthew Mathis, "Automatic TCP Buffer Tuning", ACM SIGCOMM, Oct. 1998.
- [4] Experimental TCP selective acknowledgement implementations 1998, Obtain via: <http://www.psc.edu/networking/tcp.html>
- [5] W. Richard Stevens, "TCP slow start, congestion avoidance, fast retransmit, and fast recovery algorithms", Request for Comments 2001, March 1996.
- [6] Mark Gates and Alex Warshavsky, Iperf version 1.1.1, Bandwidth Testing Tool, NLANR Applications, February 2000.
- [7] Mark Gates, "Tuning Applications for High-Performance Networks", NLANR Distributed Computing Workshop, Sep 19-21, 1999, Tucson, AZ.
- [8] S. Ostermann, M. Allman, H. Kruse, "An Application-Level Solution to TCP's Satellite Inefficiencies". Workshop on Satellite-based Information Services (WOSBIS), November, 1996. Rye, New York.
- [9] D. A. Patterson, G. A. Gibson, R. H. Katz, "The Case for Redundant Arrays of Inexpensive Disks (RAID)", Proceedings ACM SIGMOD Conference, Chicago, IL, (May 1988).
- [10] Wright Gray R., W. Richard Steven, "TCP/IP Illustrated, Volume 2 : The Implementation", Jan 1995, Addison Wesley.
- [11] S. Ostermann, "Tcptrace Version 5.2.1", Aug 1998.

ACKNOWLEDGMENTS

The authors thank Mr. Scott Wahlstrom and Mr. Marco Mazzucco for help in running certain tests for this paper.