



IEEE Global Telecommunications Conference

GLOBECOM'91

Phoenix, Arizona • December 2-5, 1991

"COUNTDOWN TO THE NEW MILLENNIUM"

Featuring a Mini-Theme on:
Personal Communications Services (PCS)

Conference Record

Vol. 3 of 3

VOLUME	DAY	SESSIONS	PAGES
1	Tuesday	1-20	1-694
2	Wednesday	21-42	695-1531
3	Thursday	43-60B	1533-2150

BEST AVAILABLE COPY



Sponsored by the
IEEE Communications Society
and the Phoenix IEEE Section

i

DEFS-ALA0010601

CAVIUM-1029
Cavium, Inc. v. Alacritech, Inc.
Page 001



BEST AVAILABLE COPY

Additional sets of Volumes 1, 2, and 3 may be ordered from:

IEEE Service Center
Publications Sales Department
445 Hoes Lane
P.O. Box 1331
Piscataway, New Jersey 08855-1331

IEEE GLOBECOM '91

IEEE Catalog No.: 91CH2980-1
ISBN Numbers: 0-87942-697-7 Softbound
0-87942-698-5 Casebound
0-87942-699-3 Microfiche

Library of Congress No.: 87-640337 Serial

COPYRIGHT AND REPRINT PERMISSIONS:

Abstracting is permitted with credit to the source. Libraries are permitted to photocopy beyond the limits of U.S. copyright law for private use of patrons those articles in this volume that carry a code at the bottom of the first page, provided the per-copy fee indicated in the code is paid through the Copyright Clearance Center, 20 Congress St., Salem, Mass. 01970. Instructors are permitted to photocopy isolated articles for noncommercial classroom use without fee. For other copying, reprint or republications permission, write to Director, Publishing Services, IEEE, 345 East 47th St., New York, NY 10017. All rights reserved. Copyright © 1991 by the Institute of Electrical and Electronics Engineers, Inc.

AN OUTBOARD PROCESSOR FOR HIGH PERFORMANCE IMPLEMENTATION OF TRANSPORT LAYER PROTOCOLS

R. ANDREW MACLEAN and SCOTT E. BARVICK¹

Bellcore, 444 Hoes Lane, Piscataway, NJ 08854

ABSTRACT: The high throughputs promised by emerging network technologies are often difficult to achieve application-to-application because of host transport protocol bottlenecks. This paper describes an experimental prototype implementation of an outboard protocol processor which eliminates these bottlenecks by performing transport layer functions in dedicated hardware. The architecture consists of separate transmit and receive CPUs, each with checksum and DMA circuits. Measurements made using an implementation of the TCP protocol indicate that this architecture can support end-to-end throughputs in excess of 11,000 packets/sec between UNIX² hosts.

1. INTRODUCTION

MEASUREMENTS of the end-to-end performance of LANs indicate that transport and network layer protocol implementations often limit throughput between communicating applications [1]. Several solutions have been proposed, the most common of which can be categorized as follows:

- 1) increase the size of the transport protocol data unit (TPDU),
- 2) optimize the protocol software,
- 3) change to a more efficient protocol and
- 4) use a more powerful host processor.

The idea behind the first method is to increase the ratio of data bits to header bits so that the protocol processing required per data bit is reduced. Depending upon the underlying network, however, this approach may lead to increased latency in the network, or the need for fragmentation and reassembly in the lower layers.

Optimization of the communications software for throughput has been discussed for the case of TCP/IP [2]. The methodology adopted is to change the existing implementation software to reduce the protocol processing overhead. Performance comparisons of optimized and non-optimized implementations of TCP/IP found in [1] indicate that significant performance increases can be realized using such techniques.

The TCP and OSI TP4 protocols have been designed primarily for robustness and utility rather than throughput. Several 'lightweight' transport layer protocols have been proposed which offer performance improvements: NETBLT [3], XTP [4], NACK [5], SNR [7], and VMTP [6] are examples of such protocols. With respect to the hardware, typically, the data link layer (to use OSI terminology) has been handled by some kind of host network adapter and the network layer and above has

been the responsibility of the host processor. Thus there are potentially two areas where processing power can be added, the host or the communications adapter.

While the use of more powerful host processors is a common solution, there is a growing trend towards increasing the front end intelligence of the I/O. There are two primary reasons for this; firstly, the I/O subsystem often demands response times which cannot be guaranteed by the host processor or would cause the host to behave inefficiently or erratically. Secondly, it is often the case that certain compute intensive functions can be performed more quickly or more cost effectively by specialized hardware than by the host processor. Several outboard processor designs have been reported [7]-[12]. Our objective for this project has been to explore the outboard approach by designing an experimental prototype processor as a platform for analyzing transport layer protocols. We call this processor the Protocol Accelerator (PA), and it is described in the sections which follow. One of our ultimate aims is to explore different high speed protocols using this processor as a platform, in order to determine the most appropriate techniques for transport of data on high speed Metropolitan Area Networks.

2. PROTOCOL ACCELERATOR

2.1 System Configuration

The Protocol Accelerator is a board on the VME system bus. Figure 1 shows how the PA integrates into the host system. On the network side, the PA is equipped with both input and output 32 bit parallel ports, each supporting data transfer rates in excess of 320 Mbits/sec. Intentionally, no media access circuitry has been included on the PA; this provides us with the capability to connect the Protocol Accelerator to a variety of network types via appropriate adapters, or, as in the case of loop-back testing, to leave out the network circuitry completely. In future communications subsystems, the transport protocol acceleration circuitry probably would physically reside on the network adaptor card.

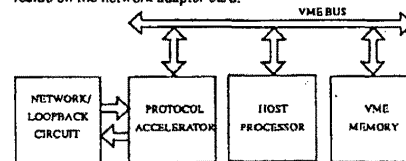


FIGURE 1. System Configuration

1. Scott Barvick is now with Wellfleet Communications, 15 Crosby Drive, Bedford, MA 01730
2. UNIX is a registered trademark of UNIX System Laboratories Inc.

49.4.1

1728

CH2980-1/91/0000-1728 \$1.00 © 1991 IEEE

GLOBECOM '91

DEFS-ALA0010603

CAVIUM-1029

Cavium, Inc. v. Alacritech, Inc.

Page 003

BEST AVAILABLE COPY

2.2 Functionality and Data Flow

The previously reported outboard processors can be categorized into single and multiple CPU architectures. The single CPU implementations [4, 10, 11] include peripheral support for high data throughputs. The multiprocessor implementations [7, 8, 9, 12] use up to eight CPUs, typically with no special peripheral devices. In our design, we exploit features from both of these approaches with a dual CPU design and special purpose peripheral circuitry in an architecture optimized for transport layer processing.

The internal functions and data flows of the protocol accelerator are shown in Figure 2. We use a dual CPU approach to protocol processing, with one CPU subsystem dedicated to the transmission, and the other to the reception. The transmit and receive CPUs are both 68020 (25 MHz) based, each with its own private resources: ROM, parallel I/O, interrupt circuitry and 128 kilobytes of random access memory (RAM). In addition there is 128 kilobytes of RAM shared by both CPUs which is also accessible to the two host busses, VME and VSB. Using both host busses simultaneously, it is possible for the PA to move data blocks both to and from host memory. The transmit and receive CPUs have VME bus master capability. All the data paths shown are 32 bits wide.

2.3 Operation

On transmit, data can be piped from one of three locations; host memory, shared memory, or transmit CPU memory into the network port, while the transmit CPU supervises transfers and compiles headers. No intermediate buffering of the application data takes place. We believe this is key to high speed operation.

On receive, parallel data from the network is pipelined to the host memory, local receive CPU memory, or shared memory. In normal operation, the receive CPU will DMA the header to local memory, perform initial processing to establish header integrity and payload destination, and then start a DMA

process to transfer the data segment of the packet either to host memory or to the shared memory region. While storage of the payload is proceeding, the receive CPU completes its processing of the header information. Messages are passed between the transmit/receive CPUs and the host either by using the shared memory region or by using interrupt mechanisms which exist among all CPUs (host, transmit or receive).

The Protocol Accelerator enables a rapid data flow to and from the network by introducing a high degree of concurrency into the communication mechanism. Several activities execute simultaneously:

- 1) host processing (of higher layers and application),
- 2) transmit protocol processing,
- 3) receive protocol processing,
- 4) data transfer from host memory to the network adapter,
- 5) data transfer from the network adapter to host memory,
- 6) receive data checksumming,
- 7) transmit data checksumming, and
- 8) MAC frame processing.

2.4 Direct Memory Access

An important feature of the hardware architecture is the dual direct memory access controllers (DMACs). The DMACs are capable of moving 32 bit data words at rates of up to 33 Mbytes/sec over VME or 30 MBytes/sec over the VSB bus directly to and from the network ports. Scatter-gather type operations are fully supported both to and from the application memory. Data paths also exist for DMA of data between the network ports and any other RAM area on the PA (i.e. CPU RAM space or shared RAM space) so that intermediate data buffering is possible whenever necessary.

2.5 Checksumming

The on-board CPUs are capable of checksumming data at a rate in the region of 75 Mbits/sec[13] but operation at this rate would leave no time for protocol processing. To maximize our data throughput, it was decided to include hardware

BEST AVAILABLE COPY

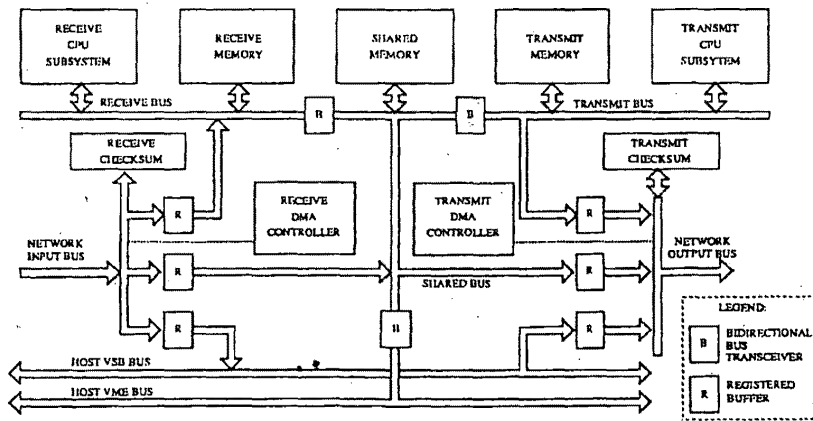


FIGURE 2. Protocol Accelerator Functional Block Diagram

checksumming in preference to using faster, more sophisticated processors for this function.

We use an 'on-the-fly' technique for the checksum and in order to take advantage of this hardware, the checksum field is required to trail the fields being checked. In the case of TCP this requires that the checksum field be moved from its usual place in the header. Although hardware schemes can be devised which leave the checksum field in place, these are somewhat more difficult to implement, and because the intent with this design was to produce a testbed suitable for many protocols, the checksumming circuitry was not designed to be protocol specific.

The circuits utilize 32 bit ones complement adders and operate in tandem with the DMA controllers; every word moved by the DMA controller is simultaneously applied to the checksum circuits. On the transmit processor, the circuit automatically inserts a 32 bit checksum at the end of both the header and the data segments. On the receiver these fields are checked, again automatically.

3. TCP IMPLEMENTATION

3.1 Introduction

The transport protocol used to demonstrate the capabilities of the Protocol Accelerator is the Transmission Control Protocol (TCP). It was chosen because it is currently one of the most wide-spread transport protocols in use on networks today. It is also the subject of a great deal of continuing research, and the expanding use of UNIX-based desktop workstations is increasing its penetration. Along with the increased use of TCP is the fear that as network rates rise, network throughput may not keep pace or may even decline as connection-oriented, window flow-controlled protocols such as TCP become a bottleneck. Some researchers do not believe that the protocol is to blame for the low throughput observed when current TCP implementations are run over experimental high speed networks [2]. Instead, they cite inefficient implementations of the protocol and interactions with the host operating system (UNIX) as the causes of the poor performance. The intensity of this debate will grow as more high performance machines running TCP observe less than ideal performance over high speed networks. Therefore, TCP was implemented to show that with an efficient implementation on the proper hardware platform, even a protocol designed for moderate-speed, error-prone networks can achieve high throughput.

3.2 Implementation Details

The high performance aspects of the Protocol Accelerator such as the dual processors, DMA, and on-the-fly checksum require the design of the transport protocol implementation to be specific to the PA. Therefore, a custom implementation of TCP was developed. The modules were written in C for the main protocol processing functions with embedded 68020 Assembly language code to perform many of the hardware specific tasks such as setting up the DMA controller, controlling timers, and polling status flags. No operating system is used on the PA. Many implementation decisions were made to optimize speed and efficiency for the 'main path' of protocol processing at the expense of processing for infrequent error conditions. These decisions are justified given the low error rates characteristic of high speed fiber networks.

The dual processor hardware architecture of the Protocol Accelerator leads naturally to the software architecture. The

TCP implementation consists of separate transmit and receive processes running their respective tasks on separate microprocessors. The transmitter and receiver do most of their processing on data stored in private memory, but they do communicate through the shared memory. The bulk of this communication occurs through the Transmission Control Block (TCB), the main TCP state information structure which resides in shared memory. At appropriate times, state changes in either the transmitter or receiver are updated in the TCB which may then be read by the other processor in the course of its work.

As noted in the hardware description, the generic nature of the PA requires that checksums be placed after the header and after the data. Other than this difference, the implementation provides all of the TCP functions required to transmit and receive data in the TCP (call) ESTABLISHED state. Among others, these functions include maintaining the retransmission queue, providing resequencing for out of order packets, supporting retransmission timers, and packetizing host data into TCP segments, or TPDUs. It must be noted that to minimize data movement, host data is moved directly between host memory and the network interface. This precludes further segmentation/reassembly of the data at what would be the Internet Protocol (IP) layer. Therefore, although certain functions of the IP layer are rolled into the TCP header (IP address, length, protocol), IP functionality is not claimed.

Another objective of this work is to quantify the effects of the end host system on outboard protocol processing. To this end a UNIX device driver, applications programming interface (API), and application were developed for the host UNIX system. The relationships among the components are shown in Figure 3.

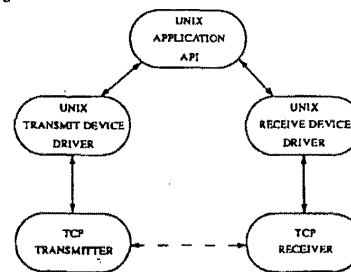


FIGURE 3. System Software Architecture

Because many variations are possible with software in the UNIX environment, attempts were made to keep the device driver, API, and application code as simple as possible while maintaining functional similarity to current methodologies for protocol/system interfaces. The results may then be used to extrapolate meaningful performance expectations of other systems with different basic parameters such as processor capability, network packet size, or bus speeds.

BEST AVAILABLE COPY

49.4.3

1730

DEFS-ALA0010605

CAVIUM-1029
Cavium, Inc. v. Alacritech, Inc.
Page 005

Explore Litigation Insights

Docket Alarm provides insights to develop a more informed litigation strategy and the peace of mind of knowing you're on top of things.

Real-Time Litigation Alerts



Keep your litigation team up-to-date with **real-time alerts** and advanced team management tools built for the enterprise, all while greatly reducing PACER spend.

Our comprehensive service means we can handle Federal, State, and Administrative courts across the country.

Advanced Docket Research



With over 230 million records, Docket Alarm's cloud-native docket research platform finds what other services can't. Coverage includes Federal, State, plus PTAB, TTAB, ITC and NLRB decisions, all in one place.

Identify arguments that have been successful in the past with full text, pinpoint searching. Link to case law cited within any court document via Fastcase.

Analytics At Your Fingertips



Learn what happened the last time a particular judge, opposing counsel or company faced cases similar to yours.

Advanced out-of-the-box PTAB and TTAB analytics are always at your fingertips.

API

Docket Alarm offers a powerful API (application programming interface) to developers that want to integrate case filings into their apps.

LAW FIRMS

Build custom dashboards for your attorneys and clients with live data direct from the court.

Automate many repetitive legal tasks like conflict checks, document management, and marketing.

FINANCIAL INSTITUTIONS

Litigation and bankruptcy checks for companies and debtors.

E-DISCOVERY AND LEGAL VENDORS

Sync your system to PACER to automate legal marketing.