

# Exploiting NIC architectural support for enhancing IP-based protocols on high-performance networks

Hyun-Wook Jin<sup>a,\*</sup>, Pavan Balaji<sup>a</sup>, Chuck Yoo<sup>b</sup>, Jin-Young Choi<sup>b</sup>, Dhableswar K. Panda<sup>a</sup>

<sup>a</sup>Computer Science and Engineering, The Ohio State University, 2015 Neil Avenue, Columbus, OH 43210, USA

<sup>b</sup>Computer Science and Engineering, Korea University, 1-5 Ka, Anam-Dong, Sungbuk-Ku, Seoul, 136-701, Republic of Korea

Received 1 May 2004; received in revised form 11 May 2005; accepted 11 May 2005

Available online 26 July 2005

## Abstract

While a number of user-level protocols have been developed to reduce the gap between the performance capabilities of the physical network and the performance actually available, their compatibility issues with the existing sockets-based applications and IP-based infrastructure has been an area of major concern. To address these compatibility issues while maintaining a high performance, a number of researchers have been looking at alternative approaches to optimize the existing traditional protocol stacks. Broadly, previous research has broken up the overheads in the traditional protocol stack into four related aspects, namely: (i) compute requirements and contention, (ii) memory contention, (iii) I/O bus contention and (iv) system resources' idle time. While previous research dealing with some of these aspects exists, to the best of our knowledge, there is no work which deals with all these issues in an integrated manner while maintaining backward compatibility with existing applications and infrastructure. In this paper, we address each of these issues, propose solutions for minimizing these overheads by exploiting the emerging architectural features provided by modern Network Interface Cards (NICs) and demonstrate the capabilities of these solutions using an implementation based on UDP/IP over Myrinet. Our experimental results show that with our implementation of UDP, termed as E-UDP, can achieve up to 94% of the theoretical maximum bandwidth. We also present a mathematical performance model which allows us to study the scalability of our approach for different system architectures and network speeds.

© 2005 Elsevier Inc. All rights reserved.

**Keywords:** Clusters; UDP/IP; Myrinet; Protocol offload; Overhead pipelining

## 1. Introduction

Commodity off-the-shelf (COTS) clusters have been accepted as a feasible and cost-effective approach to mainstream supercomputing for a broad subset of applications. Most of the success of these COTS clusters is derived from the high-performance-to-cost ratio achievable through them. With the advent of the several modern high-speed interconnects, such as Myrinet [9], InfiniBand [4], Quadrics [35],

10-Gigabit Ethernet [16,21,22] and others, the bottleneck in the data communication path in such clusters has shifted to the messaging software at the sending and the receiving side.

Researchers have been looking at alternatives by which one could increase the communication performance delivered by clusters in the form of low latency and high bandwidth user-level protocols such as the Virtual Interface Architecture (VIA) [11], FM [34] and GM [23] for Myrinet, U-Net [44] for ATM and Ethernet, EMP [40,41] for Gigabit Ethernet and others. While this approach is good for writing new applications which completely reside inside the cluster environment, these have several limitations with respect to compatibility with existing applications and infrastructure. In particular, we look

\* Corresponding author.

E-mail addresses: [jinhy@cse.ohio-state.edu](mailto:jinhy@cse.ohio-state.edu) (H.-W. Jin), [balaji@cse.ohio-state.edu](mailto:balaji@cse.ohio-state.edu), [balaji@se.ohio-state.edu](mailto:balaji@se.ohio-state.edu) (P. Balaji), [hxy@os.korea.ac.kr](mailto:hxy@os.korea.ac.kr) (C. Yoo), [choi@formal.korea.ac.kr](mailto:choi@formal.korea.ac.kr) (J.-Y. Choi), [panda@cse.ohio-state.edu](mailto:panda@cse.ohio-state.edu) (D.K. Panda).

at the following incompatibilities:

1. A number of applications have been developed in a span of several years over the traditional protocols using the sockets interface. Developing new high-performance protocols might not be directly beneficial for such applications.
2. IP is the most widely accepted and used network protocol today. However, the above mentioned user-level protocols are not compatible with existing IP infrastructures, i.e., an application using GM over Myrinet or EMP over Gigabit Ethernet cannot communicate across clusters where the intermediate nodes/switches are IP-based and do not understand these user-level protocols.
3. Several security mechanisms such as IPsec have been developed over IP. Using user-level protocols instead of IP-based protocols might require such security mechanisms to be re-developed for these new protocols.

Summarizing, several traditional applications primarily focus on portability across various platforms and typically span several clusters (sometimes known as cluster-of-cluster configurations). These applications rely on high-performance networks in order to achieve a high-performance for intra-cluster communication. At the same time they are based on IP-based protocols in order to allow inter-cluster communication. The sockets interface is a common choice for such applications. Further, applications built on other programming models such as the Message Passing Interface (MPI) [17], Distributed Shared Memory (DSM) [3,29], Global Arrays [33], etc. also have implementations utilizing the sockets layer underneath in order to achieve such compatibility.

Researchers have looked at some of these issues in a segregated manner. For example, user-level sockets over high-performance networks [6–8,30,39] and other substrates [38] have been developed to allow sockets-based applications to take advantage of the high-performance networks. This approach tries to solve the first issue (allowing existing sockets-based applications to take advantage of the high-performance networks), but does not address the remaining issues. Similarly, the Trapeze project [13,12] by Chase et al., tries to address issues two and three (compatibility with the existing IP infrastructure), but modifies the sockets interface resulting in incompatibility with existing applications. These and other related work are discussed in Section 7 in more detail.

To address these compatibility issues while maintaining a high-performance, a number of researchers have been looking at alternative approaches to optimize the existing traditional protocol stacks. Broadly, previous research has broken up the overheads in the traditional protocol stack into four related aspects, namely: (i) compute requirements and contention, (ii) memory contention, (iii) I/O bus contention and (iv) system resources' idle time.

In this paper, we first utilize the earlier proposed techniques, in particular those specific to Protocol Offload

Engines (POEs) [37,42,46], to implement a partial offload of the UDP/IP protocol stack over Myrinet to address the first two issues, i.e., compute requirements and memory contention. Next, we modify the Myrinet device driver to allow a delayed posting of descriptors in order to reduce the contention at the I/O bus. Finally, we implement a fine-grained overhead pipelining technique on the firmware of the NIC to minimize the link idle time. In this paper, we refer to this implementation as E-UDP (standing for Enhanced UDP). This IP-based protocol implementation is not only compatible with existing sockets applications, but also with the traditional UDP/IP stack.

In this paper, we focus on the following key questions:

- How does the performance of E-UDP compare with that of the traditional UDP stack?
- How does the performance of E-UDP compare with that of the existing user-level high-performance protocols?
- Would the feasibility of fine-grained pipelining in E-UDP be specific to the application communication pattern, i.e., is it formally verifiable that a fine-grained pipelining would be possible for any communication pattern?
- How does E-UDP perform for various other system and network configurations, e.g., for 10-Gigabit networks, faster I/O buses, etc?

To answer the first two questions, we analyze the performance impact of the above-mentioned techniques in UDP/IP over Myrinet. To answer the third question, we present a formal verification model and show the pipelining capabilities of the Network Interface Card (NIC) architecture in a generic communication pattern. Finally, to answer the fourth question, we propose an analytical model in order to study the performance of our design for various system and network configurations.

The remaining part of the paper is organized as follows: in Section 2 we present background information about the traditional UDP/IP implementation, POEs and the Myrinet network. In Section 3, we discuss the architectural interaction and implications of the UDP/IP protocol implementation. In Section 4, we present several solutions for the system resource contention and other inefficiencies introduced by the UDP/IP stack. We present the experimental and analytical results in Section 5, some discussion related to the broader impact of our work in Section 6, other related work in Section 7 and some concluding remarks in Section 8.

## 2. Background

In this section, we present a brief background about the traditional UDP/IP implementation, POE and the functionality of the Myrinet NIC. More details about each of these can be found in [25].

### 2.1. Traditional UDP/IP implementation

Like most networking protocol suites, the UDP/IP protocol suite is a combination of different protocols at various levels, with each layer responsible for a different facet of the communications.

To allow standard Unix I/O system calls such as `read()` and `write()` to operate with network connections, the file-system and networking facilities are integrated at the system call level. Network connections represented by sockets are accessed through a descriptor in the same way an open file is accessed through a descriptor. This allows the standard file-system calls such as `read()` and `write()`, as well as network-specific system calls such as `sendto()` and `recvfrom()`, to work with a descriptor associated with a socket.

On the transmission side, the message is copied into the socket buffer, data integrity ensured through checksum computation (to form the UDP checksum) and passed on to the underlying IP layer. The checksum computation on the sender side is usually performed during the copy operation to maximize the cache effect. The IP layer fragments the data to MTU sized chunks, constructs the IP header, and passes on the IP datagram to the device driver. After the construction of a packet header, the device driver makes a descriptor for the packet and passes the descriptor to the NIC using a Programmed I/O (PIO) operation. The NIC performs a DMA operation to move the actual data indicated by the descriptor from the socket buffer to the NIC buffer and raises an interrupt to inform the device driver that it has finished moving the data. The NIC then ships the data with the link header to the physical network.

On the receiver side, the NIC DMAs received segments to the socket buffer and raises an interrupt informing the device driver about this. The device driver hands it over to the IP layer using a software interrupt mechanism. The interrupt handler for this software interrupt is typically referred to as the bottom-half handler and has a higher priority compared to the rest of the kernel. The IP layer verifies the IP checksum and if the integrity is maintained, defragments the data segments to form the complete UDP message and hands it over to the UDP layer. The UDP layer verifies the data integrity of the message. When the application calls the `read()` operation, the data is copied from the socket buffer to the application buffer.

### 2.2. Protocol offload engines

The processing of traditional protocols such as TCP/IP and UDP/IP is accomplished by software running on the central processor, CPU or microprocessor, of the server. As network connections scale beyond Gigabit Ethernet speeds, the CPU becomes burdened with the large amount of protocol processing required. Resource-intensive memory copies, checksum computation, interrupts and reassembling of

out-of-order packets put a tremendous amount of load on the host CPU. In high-speed networks, the CPU has to dedicate more processing to handle the network traffic than to the applications it is running. POE are emerging as a solution to limit the processing required by CPUs for networking links.

The basic idea of a POE is to offload the processing of protocols from the host processor to the hardware on the adapter or in the system. A POE can be implemented with a network processor and firmware, specialized ASICs, or a combination of both. Most POE implementations available in the market concentrate on offloading the TCP and IP processing, while a few of them focus on other protocols such as UDP/IP, etc.

As a precursor to complete protocol offloading, some operating systems have started incorporating support for features to offload some compute intensive features from the host to the underlying adapters. TCP/UDP and IP checksum offload implemented in some server network adapters is an example of a simple offload. But as Ethernet speeds increased beyond 100 Mbps, the need for further protocol processing offload became a clear requirement. Some Gigabit Ethernet adapters complemented this requirement by offloading TCP/IP and UDP/IP segmentation or even whole stack on to the network adapter [1,2].

POE can be implemented in different ways depending on the end-user preference between various factors like deployment flexibility and performance. Traditionally, firmware-based solutions provided the flexibility to implement new features, while ASIC solutions provided performance but were not flexible enough to add new features. Today, there is a new breed of performance optimized ASICs utilizing multiple processing engines to provide ASIC-like performance with more deployment flexibility.

### 2.3. Myrinet network

Myrinet is a high-performance Local Area Network (LAN) developed by Myricom Incorporation [9]. In this section, we briefly describe the Myrinet NIC architecture based on LANai9. The Myrinet NIC consists of a RISC processor named LANai, memory, and three DMA engines (host DMA engine, send DMA engine, and receive DMA engine). Fig. 1 illustrates the Myrinet NIC architecture.

The LANai processor executes the Myrinet Control Program (MCP), i.e., the firmware on the NIC. The NIC memory stores the data for sending and receiving. The host DMA engine is responsible for the data movement between the host and the NIC memories through the I/O bus. On the other hand, the send DMA engine deals with moving the data from the NIC memory to the Myrinet link. Similarly, the receive DMA engine deals with moving the data from the Myrinet link to the NIC memory.

There are several emerging features provided by the Myrinet NIC. First, the programmability provided by the Myrinet NICs can be utilized to modify the implementation of existing features and/or add more features and function-

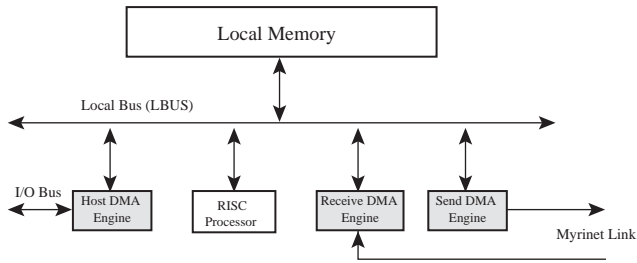


Fig. 1. Myrinet NIC Architecture.

ality to the NIC. Such programmability can provide plentiful opportunities to enhance existing IP-based protocols.

The second interesting feature is the capability of the memory on the Myrinet NIC. The memory on the NIC runs at the same clock speed as the RISC processor. Further, the LBUS (shown in Fig. 1) operates at twice the chip clock speed (two LBUS memory cycles for every clock cycle). The host DMA engine, the receive DMA engine and the send DMA engine each can request a maximum of one memory access per clock cycle. The on-chip processor can request up to two memory accesses per clock cycle. Further, the memory itself can serve up to two memory accesses per clock cycle. This means that two DMA engines, e.g., the host DMA engine and the send DMA engine, can access the memory simultaneously.

Yet another interesting feature provided by Myrinet is the capability of the host DMA engine on the Myrinet NIC. The host DMA engine allows checksum computation during the DMA operation itself. In addition, to specify the end of the buffer for a send DMA operation, the Myrinet NIC provides two kinds of registers. One is the SMLT register, which specifies not only the end of the buffer in the DMA operation but also the end of the message segment. The other register, SML, only indicates the end of the buffer. Therefore, several chunks of data sent with the SML register set, are notified as parts of the same segment on the receiver side.

### 3. Architectural viewpoint of UDP/IP

While the UDP/IP data and control path is fairly straightforward with respect to the operating system functionality, it has a number of implications on the architectural requirements of the system which implements this. These architectural requirements deal with several issues in the system such as the computational requirements of the protocol stack, memory contention caused by the stack, the I/O bus contention, etc.

#### 3.1. Interaction of UDP/IP with system resources

UDP/IP interacts with several system resources such as CPU, host memory, I/O bus and the network link. In this section, we briefly point out the requirements and extent of these interactions using the Linux UDP/IP stack as an example implementation.

*CPU interaction:* As described in Section 2.1, a number of components in the UDP/IP data path tend to be compute intensive. For example, the copy of the data from/to the application buffer occurring in the UDP/IP layer has large computation requirements. Similarly, the checksum computation occurring as a part of the bottom-half (described in Section 2.1) on the receiver side requires compute resources as well. The bottom-half typically has a higher priority compared to the rest of the kernel. This means that checksum computation for incoming packets is given a higher priority as compared to copying of the data to the application buffer. This biased prioritization of providing CPU resources for the different components has interesting implications as we will see in Section 3.2.

*Host memory interaction:* In cases where the sender buffer is touched before transmission or when the same buffer is used for transmission several times (e.g., in a micro-benchmark test), the application buffer can be expected to be in cache during the transmission operation. In such cases the data copy from the application buffer to the socket buffer is performed with cached data and does not incur any memory references. However, the case on the receiver side is quite different. On receiving data, the NIC initiates a DMA operation to move the received data into the host memory. If the host memory area for DMA is in the cache of any processor, the cache lines are invalidated and the DMA operation allowed to proceed to memory. The checksum computation, which follows the DMA operation, thus accesses data that is not cached and always requires memory accesses. Thus, the checksum computation, the DMA operation and also the data copy operation in some cases compete for memory accesses.

*I/O bus interaction:* As described in Section 2.1, once the IP layer hands over the data to the device driver, it forms a descriptor corresponding to the data and posts the descriptor to the NIC using a PIO operation over the I/O bus. The NIC on seeing this posted descriptor performs a DMA operation on the actual data from the host memory to the NIC memory. Both these operations as well as other DMA operations corresponding to incoming data use the same I/O bus and essentially contend for its ownership.

*NIC and link interaction:* The host DMA engine on the Myrinet NIC performs the DMA operation to fetch the data from the host memory to the NIC memory. The send DMA engine waits for the DMA operation to complete before it can transmit the data from the NIC memory to the link. This delay in transmitting the data can lead to the link being idle for a long period of time, in essence under-utilizing the available link bandwidth.

#### 3.2. Implications of system resources on UDP/IP performance

The interaction of the UDP/IP stack with the various system resources has several implications on the end performance it can achieve. In this section, we study these



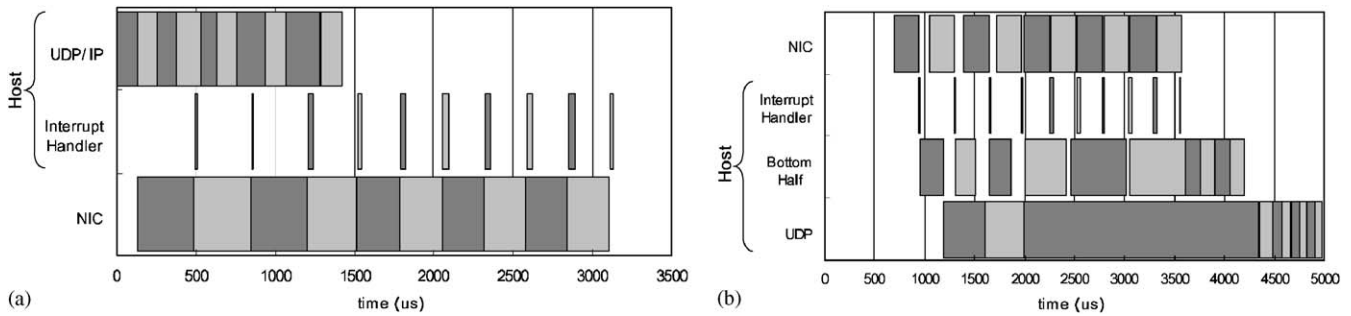


Fig. 2. Time flow chart of the host and the NIC overheads: (a) sender side and (b) receiver side.

performance implications based on some experimental results using the traditional UDP/IP stack in Linux and its time flow chart. To generate the time flow chart, we utilized the measurement methodology proposed in [27]. The key feature of this measurement methodology is to use the clock on the NIC for the overhead measurement of both the host as well as the NIC. The most important factor that should be considered for the time flow chart is that the timestamps for the host and the NIC overhead measurement have to be generated by the same clock so that we can elucidate the timing relationship between the host and the NIC overheads.

Figs. 2a and b present the time flow chart on the transmission and reception sides of the traditional UDP/IP implementation. The figures show the time flow chart of the host and the NIC when 10 UDP packets (of 32 kB each) are transmitted in a burst. The time is set to zero when the first packet is sent by the UDP application. The y-axis represents the layers that each packet goes through in order to be processed. A rectangle is the time spent in the corresponding layer to process the packet. On the receiver side, the bottom-half is not dealt as a part of UDP/IP to clearly point out the overheads caused by the checksum operation (a part of the bottom-half handler) and by the data copy operation (a part of UDP). The rectangles are shaded alternately for clarity. In these figures, we can see that the different rectangles in the same layer are of different sizes especially on the NIC on the sender side and the host on the receiver side. This is attributed to the contention between the host and the NIC for the various system resources including the host CPU, host memory, I/O bus, etc. We will deal with each of these resources in the following few subsections.

### 3.2.1. Compute requirements and contention

We first analyze the case of the host CPU contention. The host performs several compute intensive operations such as checksum computation, copying the data to the application buffer, etc. Within these, some of the operations such as the checksum computation are performed as soon as the data arrives within a higher priority context (bottom-half). During bulk data transfers where segments are received continuously, this might imply that the CPU is completely devoted to the bottom-half handler resulting in a starvation for the

other operations with compute requirements such as data copy from the socket buffer to the application buffer.

We use the results depicted in Figs. 2a and b to understand this. If we first consider the transmission side, the contention between the host and the NIC for the different system resources causes the NIC to take more time for the first four packets compared to the other packets (which is discussed in Section 3.2.3 in detail). This delay in the NIC for the first four packets is also reflected on the receiver side (Fig. 2b) where the NIC has a significant amount of idle time for the first few packets (shown as the gaps between the rectangles in the figure). These gaps in the receiver NICs active time are in turn reflected on the bottom-half handling time on the receiver side, i.e., since the packets are not coming in at the maximum speed on the network, the host can perform the bottom-half and still would have enough time to copy the data to the application buffer before the next packet comes in. In short, in this case the sender is not sending out data at the peak rate due to resource contention; this reduced data rate gives the receiver ample time to receive data and place it in the application buffer. Thus, in this case there is no contention for the CPU on the receiver side.

However, for the later packets (rectangles 5–10 for the NIC), as seen in Fig. 2a, the host has completed performing its operations; so the NIC can proceed with its operations without any contention from the host. This reduced contention for the system resources ensures that the data can be transmitted at a faster rate by the NIC, i.e., the time spent for each packet is lesser in this case. This increased transmission rate also reflects as a lesser idle time for the NIC on the receiver side (rectangles 5–10 for the NIC in Fig. 2b). Further, this reduced idle time means that the NIC continuously delivers packets to the host, thus keeping the bottom-half handler active. This results in the starvation of lower priority processes in the UDP/IP stack such as the copy of the data from the socket buffer to the application buffer. This starvation for CPU is reflected in the third rectangle in UDP of Fig. 2b where the copy operation for the third packet has to wait until all the data corresponding to the other packets has been delivered to the sockets layer.

To further verify these observations, we have also used the Performance Measurement Counters (PMCs) for the

Pentium processor to measure the impact of CPU starvation, cache miss effects, etc. on the actual time taken by the data copy operation. However, due to space restrictions, we do not present the results here and refer the reader to [25] for the same.

### 3.2.2. Memory contention

We next look at the host memory contention. Several operations in the UDP/IP stack such as checksum computation, data copy to the application buffer, etc., as well as the DMA operations to and from the network adapter compete for memory accesses. We again refer to Fig. 2 for understanding these contention issues. Observing Fig. 2b, we notice that when the rate of the incoming packets increases, the time taken for the bottom-half increases (rectangles 4–6 in the bottom-half). This is because the checksum computation in the bottom-half handler competes for memory accesses with the DMA operations carried out by the NIC (rectangles 5–10 for the NIC overhead in Fig. 2b). This figure shows the impact memory contention can have on the performance of the UDP/IP stack.

Again, to re-verify these results, we have used PMCs to measure the actual checksum computation overhead and the wait time for fetching data from the memory to cache. The results for the same can be found in [25].

### 3.2.3. I/O bus contention

Posting of descriptors by the device driver as well as DMA of data by the NIC to/from the NIC memory uses the I/O bus causing contention. The transmission side in Fig. 2a shows the increased time taken by the NIC for rectangles 1–4. However, this increase in the time taken by the NIC can be because of both I/O bus contention and memory contention. In order to understand the actual impact of the I/O bus contention, we modified the UDP/IP stack to offload the checksum computation to the NIC and allow a zero-copy implementation. These modifications completely get rid of the memory accesses by the host avoiding any memory contention that might be possible.

Fig. 3 shows the I/O bus contention for the modified UDP/IP stack. We can see that the UDP/IP overhead is negligible in this case because of the offloading of the data touching components. On the other hand, the NIC overhead for the first rectangle is significantly larger than the rest of the rectangles due to the contention between the posting of the descriptors and the DMA of the actual data to the NIC. Since the host does not touch its memory at all, we can say that this overhead is completely due to the I/O bus contention.

### 3.2.4. Link idle time

Fig. 4 shows the time flow chart for the NIC and the link overheads for the traditional UDP/IP stack. As discussed earlier, the earlier rectangles (1–4) showing the NIC overhead on the transmission side are larger than the later ones (5–10) because of the memory and I/O bus contentions. This

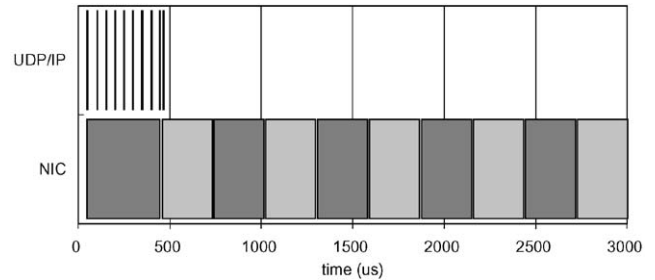


Fig. 3. Time flow chart of the host and the NIC overheads on the transmission side with the modified UDP/IP stack (checksum offloaded and zero-copy implementation).

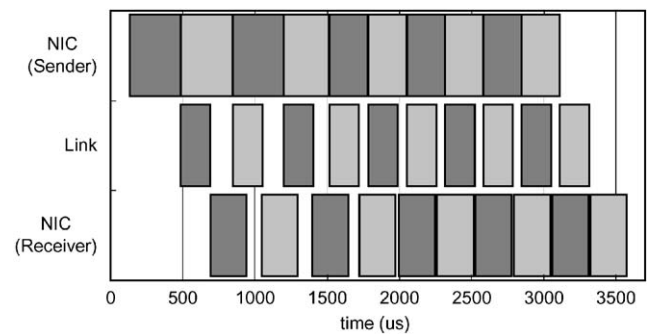


Fig. 4. Time flow chart for the NIC and the link overheads for traditional UDP/IP.

overhead on the sender NIC is reflected as idle time for the link and the receiver NIC (shown as gaps between rectangles in the figure). This figure shows that memory and I/O contention at the host can result in the link being idle for nearly 50% of the time, in essence dropping the effective link bandwidth to half. Further, the startup overhead for the link (shown as the gap before the first rectangle) becomes especially important for applications which do not carry out bulk data transfers, e.g., latency micro-benchmark test.

Overall, these results attempt to demonstrate the various inefficiencies present in the current UDP/IP implementation due to contention for the existing system resources. In Section 4, we present several solutions to address these issues and implement a high-performance UDP/IP stack using the features provided by the Myrinet network.

## 4. Enhancing UDP/IP performance: Design overview

In the following few subsections, we discuss several design solutions for the issues pointed out in the previous section. In particular, we look at the issues of (i) compute requirements and contention, (ii) memory contention, (iii) I/O bus contention and (iv) system resources' idle time and suggest efficient designs in the context of the Myrinet network to address these issues.

#### 4.1. Partial protocol offloading for avoiding CPU and memory contention

As discussed in Section 3, the host CPU performs per-byte operations such as data copy and checksum computation which result in not only an intensive consumption of CPU resources, but also memory contention with the DMA operations that access the host memory from the NIC. Further, as described earlier, checksum offload and zero-copy transmission and reception implemented by some POEs allows a significant reduction in these resource usage requirements. We regard basic checksum offloading and zero-copy data transmission as pre-requisites to our work and implement these over Myrinet.

We can consider two design alternatives to offload the checksum computation on to the NIC. The first alternative is to let the NIC firmware perform the checksum computation on the data in the NIC memory. This can be performed in parallel with the protocol processing by the host CPU. Though this approach would reduce the contention for the CPU requirement on the host, it would merely shift the memory contention issue from the host memory to the NIC memory causing the checksum computation to contend with the host and network DMA operations. In addition, since the NIC processor is significantly slower than the host CPU, it is difficult to expect this approach to reduce the overhead for checksum computation.

The other approach is to utilize the checksum computation capability of the host-DMA engine. The checksum computation by the DMA engine does not generate any additional computation overhead or memory accesses since the data checksum is performed together with the DMA operation and before being placed into the NIC memory. This allows us to achieve both offloading as well as high performance. At the same time, the checksum offloading does not introduce any additional memory contention on the NIC memory.

Together with checksum computation, the data copy present in the traditional UDP implementation is another source of CPU and host memory contention. One way to remove the copy operation is to utilize the memory mapping mechanism between the application and the kernel memory spaces and use the mapped kernel memory space for all data communication [14]. Another alternative is to directly move the data in the application buffer through DMA with the address and the length of the application buffer [47]. Without the copy operation, the protocol layer can simply take the virtual address of the application buffer and translate it into the corresponding physical address by traversing the page directory and table. This physical address is used by the NIC to perform the DMA operations. In either case, since we allow applications to use arbitrary application buffers, we have to consider the linearity of the buffer. In cases where the buffer is not physically linear, we construct gather or scatter lists so that the NIC can perform a DMA operation for each linear chunk in the application buffer separately. Both approaches have their own

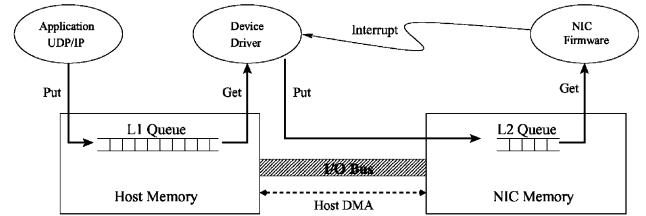


Fig. 5. Two level queues to control I/O bus contention.

advantages and disadvantages, but for efficiency and ease of implementation, we chose the second approach.

#### 4.2. Scheduling I/O bus transactions

As mention in Section 2.1, once the IP layer presents the information about the outgoing data to the device driver, it forms a descriptor for the data message and pushes this descriptor to the NIC through a PIO operation. While this approach is quite efficient when there is little contention on the I/O bus, this might lead to a significant degradation of the performance when the I/O bus gets more heavily utilized. For example, bursty communication requests from applications can generate a high rate of descriptor posting to the network adapter with each descriptor posting requiring a PIO transaction across the I/O bus.

In order to resolve this problem, we propose a two-level queue structure to allow an efficient scheduling of the rate of descriptor posting without sacrificing the performance. Fig. 5 shows the basic structure of two-level queue, where the level-1 (L1) queue is large and located in the host memory, while the level-2 (L2) queue is small and located in the NIC memory.

When an application issues a send request, if the L2 send queue is not full, the device driver puts the descriptor corresponding to the request in the L2 send queue through a PIO operation. On the other hand, if the L2 send queue is full, the descriptor is just queued in the L1 queue. When the NIC firmware completes a send, an interrupt is triggered, which means that a vacancy is created in the L2 send queue. Accordingly, the interrupt handler gets a descriptor from the L1 send queue and puts it into the L2 send queue through a PIO operation. This approach ensures that the L1 send queue behaves in a self-clocking manner, i.e., it uses the interrupt for send completion as a clock to put a new descriptor into the L2 send queue. The description of the two-level queue on the receiver side is similar.

Consequently, the two-level queue structure can mitigate the rate of PIO to less than

$$Rate_{\max} = \frac{n + \text{sizeof}(L2)}{t},$$

where  $n$  is the maximum number of requests that the NIC can process in time  $t$ . And  $\text{sizeof}(L2)$  is the number of requests that L2 queue can hold.

The kernel and the NIC do not need to know any details about the queue structure because the device driver allocates and manages the L1 queue and can decide the size of the L2 queue. Accordingly, it is easy to apply to various operating systems and NICs. While we understand the importance of the performance impacts the size of the L2 queue might have, in this paper, we do not deal with its variation and fix it to a predefined value for all experiments.

#### 4.3. Fine-grained pipelining to minimize link idle time

As mentioned earlier, the Myrinet NIC memory allows two memory accesses per clock cycle while each DMA engine can only request one memory access per clock cycle. This means that the transmission of a data segment by the send DMA engine can be pipelined with the DMA of the next segment by the host DMA engine (since these engines can read and write to the memory simultaneously).<sup>1</sup> The current GM architecture follows this approach for data transmission. We term this form of pipelining as *coarse-grained pipelining*. Earlier generation protocols over Myrinet such as Berkeley-VIA [11] achieve a much lower performance as they do not utilize this information about the memory clock speed at all and carry out the host DMA and the network DMA in a serialized manner [26]. However, as discussed in Section 3, even a coarse-grained pipelining approach does not fully avoid the idle time in the link since the transmission is not carried out till the DMA of the entire segment is completed.

In order to address the issue of the link idle time, we propose a fine-grained pipelining approach between the I/O bus and the Myrinet link. In the fine-grained pipelining approach, the NIC initiates its DMA operation as soon as a sufficient number of bytes have arrived from the host memory (through a host DMA) or the network (through a receive DMA). For instance, on the sender side the send DMA engine can send a packet out to the network while the host DMA engine is still moving a later part of the same packet from the host to the NIC memory. This approach allows us to fully pipeline the overheads of the send DMA and the host DMA for the same packet. It is to be noted that this approach only aims at reducing the per-byte overhead of the packet and does not impact the per-packet overhead associated with UDP/IP.

In order to achieve a fine-grained pipelining between the NIC and the link overheads, the following conditions should be satisfied. First, the NIC architecture has to allow multiple accesses to the NIC memory at a time by different DMA engines. The Myrinet NIC allows two DMA engines to

access the NIC memory simultaneously. Second, the NIC has to provide an approach to send several fragments of a packet separately and realize that these fragments belong to the same packet on the receiver side. As described in Section 2, the Myrinet NIC provides two kinds of registers, namely SMLT and SML to achieve this. Third, the NIC firmware should force the host DMA engine to perform the DMA operation in parallel with the network DMA engines within the same packet.

While it is easily verifiable that the NIC satisfies the first two conditions, it is very difficult to verify that the NIC firmware always guarantees the third property irrespective of the communication pattern followed by the application. To address this, we propose a formal verification for the fine-grained pipelining-based firmware model on the Myrinet NIC firmware, MCP. The MCP performs coarse-grained overhead pipelining for IP-based protocols, where DMA overheads across different packets are overlapped.

The MCP consists of multiple threads: SDMA, RDMA, SEND, and RECV. The SDMA and SEND threads are responsible for sending. The SDMA thread moves data from the host memory to the NIC memory, and the SEND thread sends data in the NIC memory to the physical network. The receiving of data is performed by the RDMA and RECV threads. The RECV thread receives data from the physical network to the NIC memory. The RDMA thread moves the received data to the host memory.

Based on this design we suggest an extended firmware design for fine-grained pipelining as shown in Fig. 6. In the figure, the states in the rectangle (with the dotted line) are the newly defined states for fine-grained overhead pipelining. The shaded states in the figure are dispatching entries of each thread. Each thread starts from the initial state and when it reaches a dispatching state, yields the processor to another thread that is ready to run without waiting for the next event to translate the state of the running thread. The yielding thread starts at a later point from the state in which the thread was suspended right before.

For example, let us consider the states for fine-grained pipelining on the RDMA thread. The initial state of the RDMA thread is the `Idle` state. If the amount of data arrived is more than the threshold for fine-grained pipelining, the state is moved to `Fine_Grained_Rdma`, where the RDMA thread initiates the host DMA operation to move the data in the NIC memory to the host memory. After finishing this DMA operation, in the `Fine_Grained_Rdma_Done` state, if there is still more data than the threshold, the RDMA thread performs the host DMA operation again moving to the `Fine_Grained_Rdma` state. Otherwise, if the receive DMA has completed, the state of the RDMA thread is changed to the `Rdma_Last_Fragment` state. In this state, the RDMA thread does the DMA operation for the rest of the packet regardless of its size.

In order to verify the correctness of the proposed state transition diagram, we used the Spin verification tool. Spin [20] is a tool for formal verification of distributed software

<sup>1</sup> In theory, unlike dual ported memory interfaces, the Myrinet NIC memory does not provide truly concurrent access. However, the number of requests the memory can handle is twice the number of requests each DMA engine can generate. So, the accesses can be assumed to be concurrent.



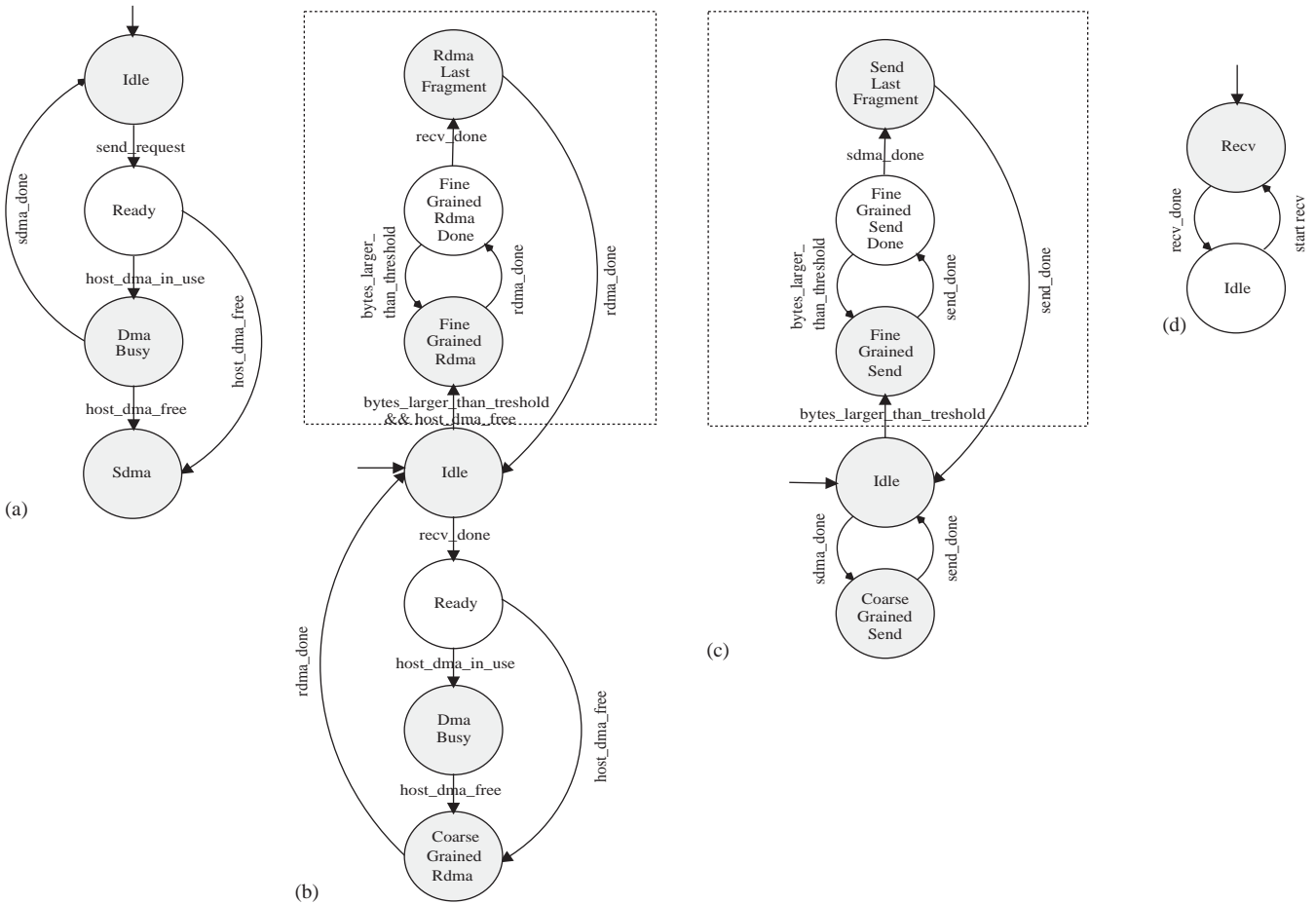


Fig. 6. State transition model for: (a) SDMA, (b) RDMA, (c) SEND, and (d) RECV threads.

systems, such as operating systems and data communication protocols. It uses a high level language to specify a system, called PROMELA (PROcess MEta LAnguage) [19]. Given a model specified in PROMELA, Spin verifies that the model satisfies properties written in linear temporal logic (LTL) [32]. LTL was designed for expressing temporal ordering of events and variables; it characterizes the behavior of systems by specifying some behavior that all traces of a system must obey.

For the formal verification, we first translate the state transition diagrams of Fig. 6 into specifications written in PROMELA. We then define propositional symbols and derive verification formulas written in LTL as Table 1.

Formulas 1 and 2 represent the properties that the suggested model performs the fine-grained pipelining. Formula 3 represents that the model utilizes the Myrinet link in full-duplex mode. Formula 4 ensures correctness, i.e., only one thread between SDMA and RDMA should occupy the host DMA engine at a time. Using these formulas with Spin, we formally verified that the above presented model performs fine-grained pipelining with any generic communication pattern.

#### 4.4. Performance modeling

To analyze the performance of our design described in Sections 4.1–4.3 on various systems, we propose a mathematical performance model. First, we derive the performance model for coarse-grained pipelining as a base model to compare against. Next, we describe the performance model for fine-grained pipelining. Both models implement the partial protocol offloading as well as the two-level queuing and differ only in the pipelining mechanism. In this section, we model our implementation based on the Myrinet network architecture. However, as we will see in Section 6, many of these features are quite general and supported by several emerging NICs for high-speed networks. Therefore, the performance models in this section are expected to give us a strong hints about the benefits achievable by our proposed protocol offloading and overhead pipelining mechanisms on next generation high-speed networks.

##### 4.4.1. Coarse-grained pipelining

In the coarse-grained pipelining model, pipelining between the overheads for the  $(p + 1)$ th packet at the  $i$ th layer and the  $p$ th packet at the  $(i + 1)$ th layer occurs as shown in

Table 1  
Linear temporal logic formulas

Propositional symbols	<pre>#define sdma (SDMA_state == Sdma) #define coarse_rdma (RDMA_state == Coarse_Grained_Rdma) #define fine_sdma (RDMA_state == Fine_Grained_Rdma) #define coarse_send (SEND_state == Coarse_Grained_Send) #define fine_send (SEND_state == Fine_Grained_Send) #define recv (RECV_state == Recv)</pre>
Formula 1	<pre>&lt;&gt; ( sdma &amp;&amp; ( coarse_send    fine_send ) ) Can the SEND thread initiate a send DMA while the SDMA thread performs a host DMA, and vice versa?</pre>
Formula 2	<pre>&lt;&gt; ( ( coarse_rdma    fine_rdma ) &amp;&amp; recv ) Can the RECV thread initiate a receive DMA while the RDMA thread performs a host DMA, and vice versa?</pre>
Formula 3	<pre>&lt;&gt; ( ( coarse_send    fine_send ) &amp;&amp; recv ) Can the SEND thread initiate a send DMA while the RECV thread performs a receive DMA, and vice versa?</pre>
Formula 4	<pre>[ ] ( sdma &amp;&amp; ( coarse_rdma    fine_rdma ) ) The SDMA thread cannot use the host DMA engine while RDMA thread utilizes it, and vice versa.</pre>

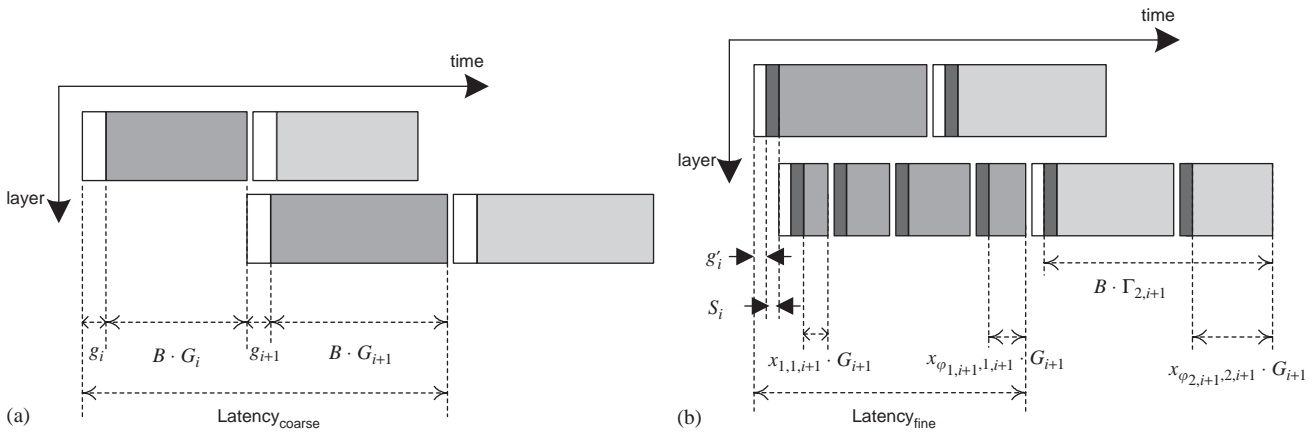


Fig. 7. Overhead pipelining: (a) coarse grained and (b) fine grained.

Fig. 7a, where the smaller numbered layer is the upper layer. In the figure,  $g_i$  and  $G_i$  denote the per-packet and the per-byte overheads at the  $i$ th layer, respectively.  $B$  is the byte size of the packet.

In this case, the one-way latency per packet is given by

$$Latency_{coarse} = \sum_{i=1}^n (g_i + B \cdot G_i) \quad (1)$$

and the bandwidth is given by

$$\begin{aligned}
 Bandwidth_{coarse} &= \lim_{m \rightarrow \infty} \frac{B \cdot m}{(m-1) \cdot (g_b + B \cdot G_b) + \sum_{i=1}^n (g_i + B \cdot G_i)} \\
 &= \frac{B}{g_b + B \cdot G_b}, \quad (2)
 \end{aligned}$$

where  $n$  is the number of layers that perform the overhead pipelining. The subscript  $b$  represents the bottleneck layer, and  $m$  is the number of packets. We analyze the bandwidth

for a large number of packets (i.e.,  $m \rightarrow \infty$ ) since we assume that the test would use a massive data transmission.

#### 4.4.2. Fine-grained pipelining

Based on the implementation of partial protocol offloading and two-level queuing, we can achieve fine-grained pipelining with the firmware model described in Section 4.3. In the fine-grained pipelining approach, a layer initiates its per-byte processing as soon as a sufficient number of bytes have arrived from the upper layer. Therefore, the overhead of a packet at the  $i$ th layer and the same packet at the  $(i+1)$ th layer are fully pipelined except the per-packet overhead as shown in Fig. 7b. In the figure,  $S_i$  is the start-up overhead that is required to start a per-byte operation such as a DMA operation,  $g'_i$  is the per-packet overhead excepting the start-up time,  $S_i$ , (i.e.,  $g'_i = g_i - S_i$ ), and  $x_{f,p,i}$  is the size of  $f$ th fragment of the  $p$ th packet in the  $i$ th layer. A characteristic of fine-grained pipelining is that the per-byte processing of a layer is affected by the per-byte overhead of the upper layer. Accordingly, we define a new parameter,  $\Gamma_{p,i}$  [25], which represents the per-byte processing time.

Then, the one-way latency of the fine-grained pipelining is as follows:

$$Latency_{\text{fine}} = \sum_{i=1, i \neq e}^n (g_i + t \cdot G_i) + g'_e + B \cdot \Gamma_{1,e}, \quad (3)$$

where  $e$  is the last layer such that  $\Gamma_{1,e} \neq 0$ , and  $t$  is the threshold that defines the minimum number of bytes making a layer start its per-byte processing.

In order to model the bandwidth, we define a parameter,  $u$ , which is the number of packets after which the fine-grained pipelining is changed into the coarse-grained pipelining. This occurs because when a layer passes the data to the next layer which happens to be slower, data accumulates in the buffer of the next layer triggering progressively larger transfers. Therefore, for all packets from the first to the  $u$ th, the per-byte overhead is defined by  $\Gamma_{p,i}$  while, after the  $u$ th packet, the per-byte overhead is the same as that of coarse-grained pipelining (i.e.,  $G_i$ ). Thus, the bandwidth of fine-grained pipelining is given by

$$\begin{aligned} Bandwidth_{\text{fine}} &= \lim_{m \rightarrow \infty} \frac{B \cdot m}{\sum_{i=1}^{n-1} g'_i + \sum_{p=1}^m (g'_n + B \cdot \Gamma_{p,n})} \\ &\approx \lim_{m \rightarrow \infty} \frac{B \cdot m}{\sum_{i=1}^{n-1} g'_i + \sum_{p=1}^m (g'_b + B \cdot G_b + k_{p,b} \cdot S_b)} \\ &= \lim_{m \rightarrow \infty} \frac{B \cdot m}{\sum_{i=1}^{n-1} g'_i + m \cdot (g'_b + B \cdot G_b) + \sum_{p=1}^u (k_{p,b} \cdot S_b) + \sum_{p=u}^m S_b} \\ &= \frac{B}{g_b + B \cdot G_b}, \end{aligned} \quad (4)$$

where  $k_{p,i}$  is the number of fragments of the  $p$ th packet of the  $i$ th layer [25].

An interesting result is that Eq. (4) is the same with Eq. (2). This is because fine-grained pipelining switches to coarse-grained pipelining after the  $u$ th packet when a large number of packets are transmitted in burst. As a result, fine-grained pipelining can achieve a low latency without sacrificing the bandwidth.

More detailed proofs, equations and explanations for fine-grained pipelining are skipped in this paper in order to maintain simplicity and are presented in [25].

## 5. Experimental results

In this section, we first present the capability of our E-UDP implementation to avoid the various inefficiencies pointed out in Section 3. Next, we present the performance achieved by E-UDP compared to other protocol implementations. Finally, we present the results of our analytical model showing the performance of E-UDP for various system and network configurations.

For the evaluation, we used a pair of machines equipped with an Intel Pentium III 1 GHz processor on an ASUS motherboard (Intel 815EP chipset). Each machine has a

Myrinet NIC (LANai 9.0) set to a 32 bit 33 MHz PCI slot, and the NICs are directly connected to each other through the Myrinet-1280 link. The Linux kernel version used is 2.2, and we adopt GM (version 1.4) for the device driver and the firmware on the Myrinet NIC. The MTU size is set to 32 kB.

### 5.1. System resource consumption in E-UDP

To analyze the effect of partial protocol offloading on the CPU overhead, we measured the overhead on both the host and the NIC CPUs. Figs. 8a and b compare the CPU overheads of the original UDP and E-UDP for small (1 B) and large (32 kB) message sizes, respectively. For small messages, though the copy operation and the checksum computation overheads are small, we can see a slight reduction in the CPU overhead, especially on the receiver side. On the other hand, the NIC overheads of E-UDP are increased due to the offload of the copy operation and the checksum computation. Further, some other functionalities such as a part of the UDP/IP header manipulation also have been moved to NIC. Overall, the accumulated overhead on both the host and the NIC are nearly equal (44.7us on E-UDP vs. 42.3us on original UDP).

For large messages, however, we can see a large benefit through partial protocol offloading. By offloading the per-byte operations from the host, we achieve a very small host overhead regardless of the message size. At the same time, there is no significant increase in the overhead on the NIC.

To observe whether E-UDP resolves the resource contention and the idle resource problems, we study the time flow chart of E-UDP for 10 packets each of 32 kB size. The time flow chart of the original UDP has already been shown in Section 3. Fig. 9 shows the time flow chart of E-UDP. We can see that the overhead of each layer can be fully pipelined with the others from the sender side to the receiver side. This is due to the fact that E-UDP eliminates the CPU, the memory, and the I/O bus contentions. In addition, E-UDP performs a fine-grained overhead pipelining to overlap the NIC and the link overheads. Consequently, the largest overhead (i.e., the NIC overhead) hides the smaller overheads and allows us to achieve a performance close to the theoretical maximum.

### 5.2. Latency and bandwidth

In this section, we compare the performance of E-UDP with that of the original UDP, GM and Berkeley-VIA [11], a well-known implementation of VIA. Since there is no implementation of Berkeley-VIA on LANai9, we measured its performance with a LANai4-based Myrinet NIC on the same platform.

Fig. 10 compares the latency of E-UDP, GM, Berkeley-VIA, and original UDP with the theoretical minimum latency of the experimental system for large and small messages, respectively. The latency test is conducted in a

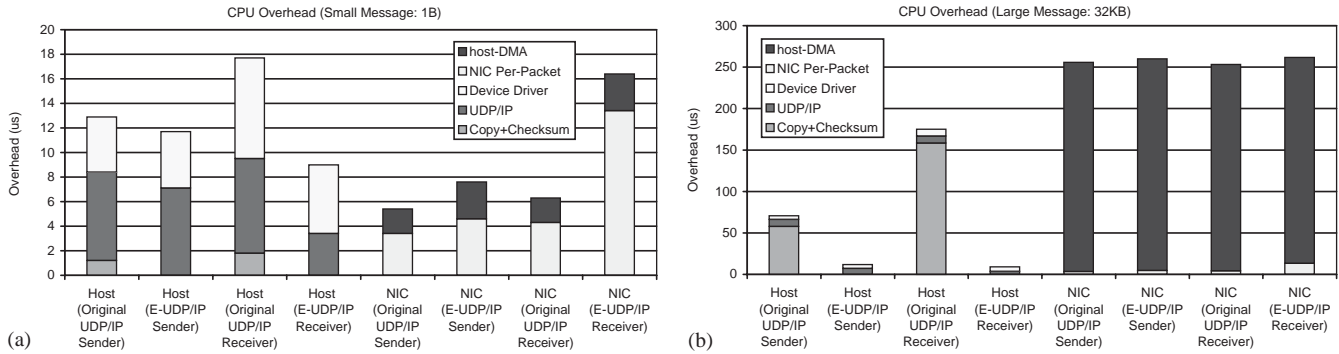


Fig. 8. CPU overhead on the host and the NIC: (a) small messages (1 byte) and (b) large messages (32 kbytes).

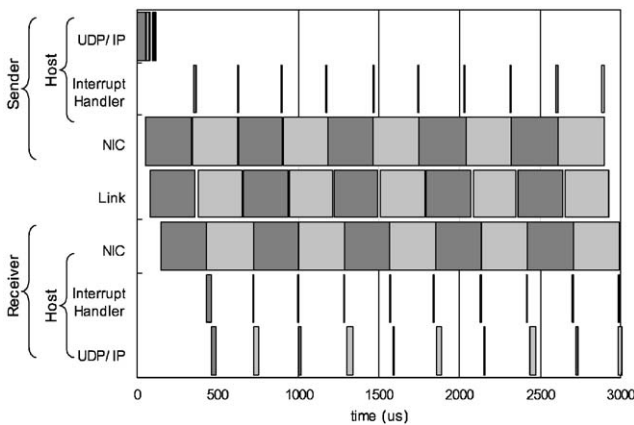


Fig. 9. Time flow chart for E-UDP.

ping-pong fashion, and the results are derived from the round-trip time by dividing it by two. Since the bottleneck of the system is the PCI bus (1007 Mbps of PCI vs. 1280 Mbps of Myrinet-1280), the theoretical maximum performance in this subsection is derived from that achievable by the PCI bus used.

We can observe that the latency of E-UDP is smaller than the others and almost even with the theoretical minimum latency. An interesting result is that the latency of E-UDP is even smaller than that of the user-level protocols, such as GM and Berkeley-VIA for large message sizes. This is because E-UDP performs fine-grained overhead pipelining on the NIC. On the other hand, in the case of small data sizes (Fig. 10b), GM shows the smallest one-way latency. This is because for small messages, the per-packet overhead becomes the dominant factor. User-level protocols have a lower per-packet overhead compared to UDP/IP following which they are able to achieve a lower latency.

Fig. 11 shows the bandwidth achieved by E-UDP compared to the other protocols. We have used the *ttcp* [43] and the *gm\_allsize* benchmarks to measure the bandwidth of E-UDP and GM, respectively. We measured the throughput of Berkeley-VIA using the *window* program. *gm\_allsize* is a benchmark provided as a part of the Myrinet software

package to measure the performance achievable by GM. *window* is a test program included in Berkeley-VIA software package for Linux. We turned on the `LAZY_BENCHMARK` option for the *window* program; this generates the theoretical maximum bandwidth achievable by evaluating the maximum injection rate to the network on the sender side regardless of the receiver.

A notable result is that E-UDP achieves a peak bandwidth of 951 Mbps which is about 94% of the theoretical maximum bandwidth of the experimental system. The bandwidth of GM is slightly lower than that of E-UDP. This is because GM splits a large packet into several segments with a fixed size of 4 kB. This segmentation increases the packet processing overhead in proportion to the number of segments. Consequently, the throughput suffers.

An unexpected result in Fig. 11 is that for large data sizes (larger than 1 kB) Berkeley-VIA shows a much lower throughput than the others. This can be partly because we measured its performance with a LANai4-based NIC. However, since the rest of the experimental system was the same, the data-touching portion is unchanged between LANai4 and LANai9. Based on this, the drop in bandwidth is attributed to the lack of pipelining between the NIC and the link overheads, i.e., the NIC firmware of Berkeley-VIA serializes DMA operations even for those performed by different DMA engines without taking advantage of the capabilities of the Myrinet NIC memory to allow simultaneous access to both the DMA engines.

### 5.3. Performance modeling results

To verify that the mathematical performance model suggested in Section 4.4 is accurate, we compare the performance evaluated from the model with the real performance numbers on our test-bed. Figs. 12a and b show the latency and bandwidth comparisons, respectively. Since, from the performance model equations in Section 4.4, the bandwidth of the coarse and the fine-grained pipelining mechanisms are the same, Fig. 12b does not deal with coarse and fine-grain pipelining separately. Real data for coarse-grain pipelining refers to a version of UDP with partial protocol offloading



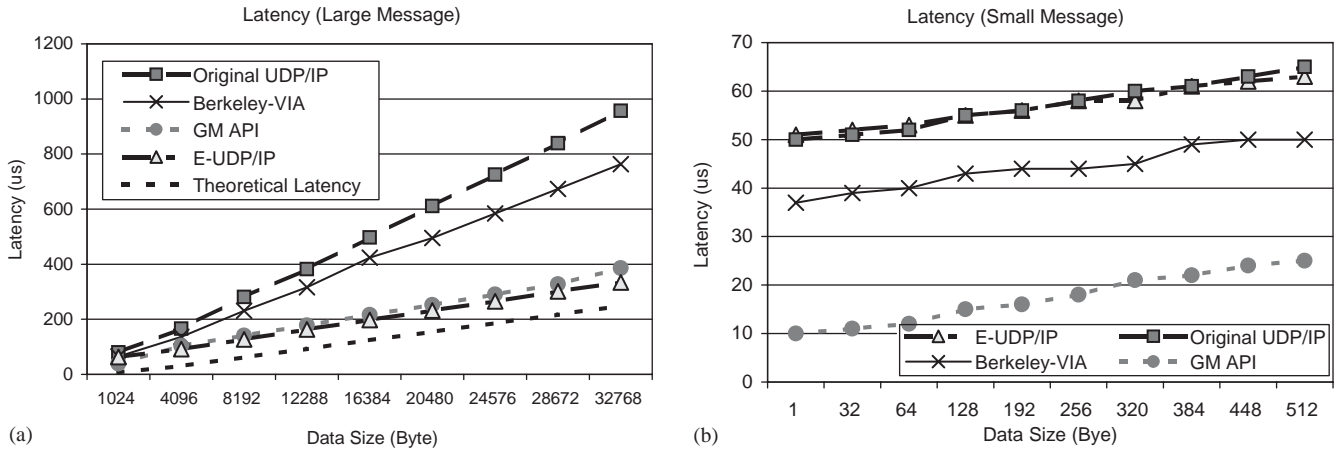


Fig. 10. Latency comparison for (a) large messages and (b) small messages.

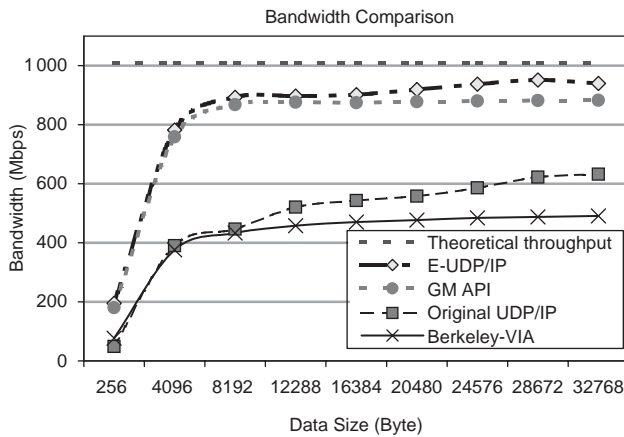


Fig. 11. Bandwidth comparison.

and two-level queuing for the I/O requests, i.e., it differs from E-UDP in only the pipelining mechanism. As we can see in the figures, the performance model matches the actual data very closely with an error of less than 5%.

Based on the performance model, we tried to analyze the latency and bandwidth of E-UDP on faster networks and I/O buses than our evaluation system. We considered 2 and 10 Gbps networks to reflect the characteristics of the emerging networks such as Myrinet-2000, 10-Gigabit Ethernet, and InfiniBand. In addition, we took account of 64 bit/66 MHz and 64 bit/133 MHz PCI systems. We used the per-packet overhead values measured on E-UDP and decided the per-byte overhead values according to the target network and the PCI bus speed.

Figs. 13a and b show the latency and bandwidth on a 2 Gbps network with a 64 bit/66 MHz PCI bus. We can see that in this case the fine-grained pipelining can achieve a very low latency for large message sizes compared to coarse-grained pipelining. In addition, the rate of increase of the overhead is equal to that of the network link latency; this shows that fine-grained pipelining is able to hide the DMA

overhead on the NIC behind the link overhead successfully. Moreover, Fig. 13b shows that both overhead pipelining mechanisms can achieve a near physical bandwidth.

The results on a 10 Gbps network and a 64bit/133 MHz PCI bus are shown in Figs. 14a and b. We can see that even fine-grained pipelining cannot fully utilize the network link. The main reason for this is that the network and I/O bus systems are an order of magnitude faster than the processor on the NIC (the modeling is based on the Myrinet LANai9 processor); this results in a relatively high percentage of the overall overhead to be associated with the protocol processing time (per-packet overhead). With the current faster processors on the NIC and hardware-based solutions being incorporated for protocol processing, we expect this overhead to be reduced significantly.

## 6. Discussion

While this work has been implemented and evaluated with the Myrinet network architecture, several of these ideas are quite general and applicable in various other architectures. In this context, we would like to recognize and consolidate a primary set of features provided by the Myrinet and other similar architectures which we feel can form a basis for a generic POE. In particular, we would like to discuss: (i) integrated checksum and DMA capability, (ii) exposing Scatter/Gather capabilities to the device driver, (iii) solicit bit for message segmentation, (iv) multiple DMA engine support and (v) memory bandwidth for internal data movement on the network adapter.

*Integrated Checksum and DMA capability:* As mentioned in Section 5, the performance of a protocol stack for small messages mainly depends on the per-packet overheads associated including the protocol stack data path overhead, interrupts, etc. On the other hand, the performance for large messages depends on the per-byte overhead associated including checksum, copy and data DMA.

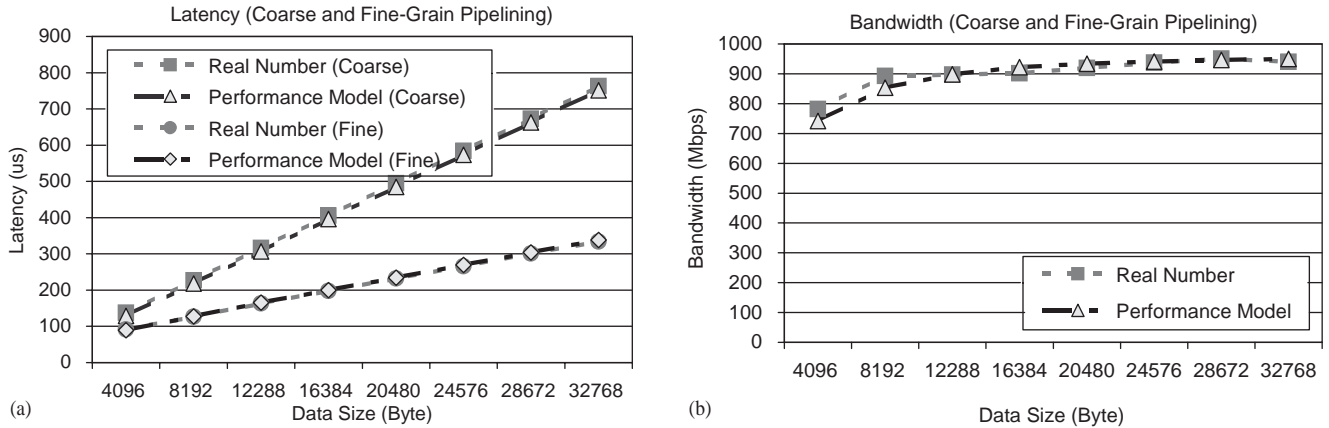


Fig. 12. Performance modeling results verification: (a) latency and (b) bandwidth.

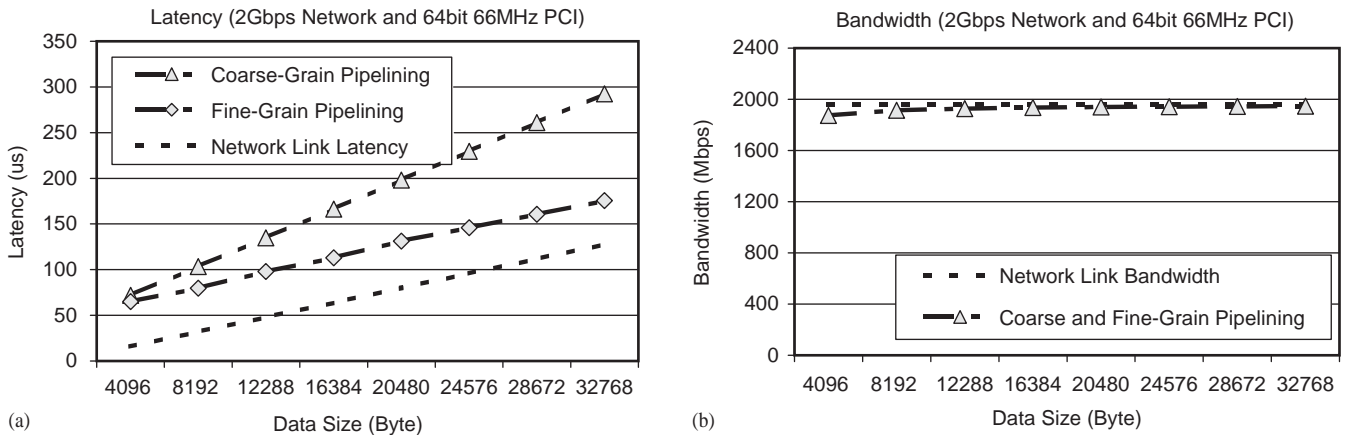


Fig. 13. Latency and bandwidth measurements on a 2 Gbps network with a 64 bit/66 MHz I/O bus.

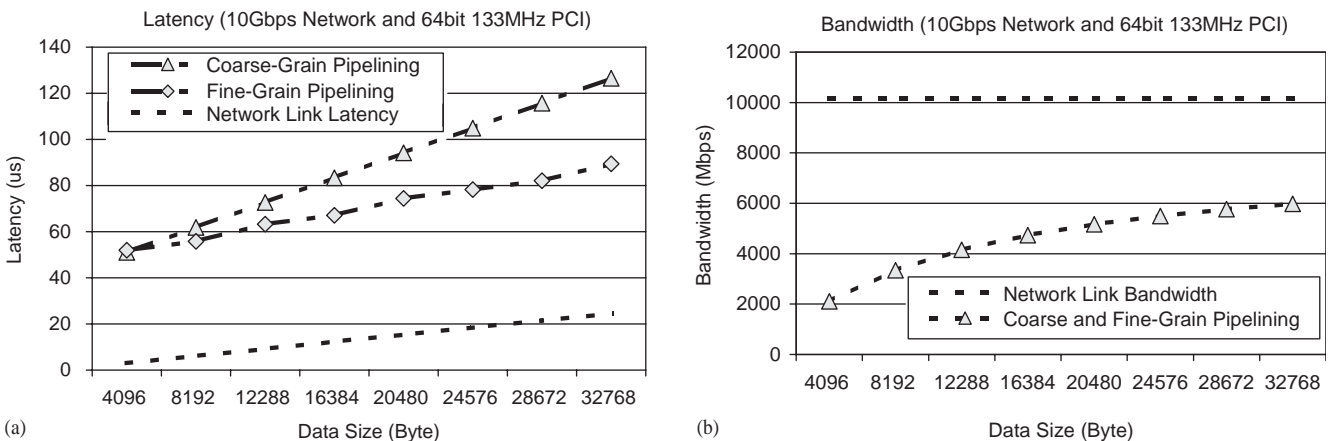


Fig. 14. Latency and bandwidth measurements on a 10 Gbps network with a 64 bit/133 MHz I/O bus.

Some network adapters (e.g., some of the Gigabit Ethernet adapters available in market) allow offloading the checksum computation to the NIC. Though this technique is quite beneficial for the host in terms of the CPU required for checksum processing, this solution merely shifts this

overhead from the host CPU to the network adapter. Further, this requires additional fetching and manipulation of data on the memory in the network adapter, increasing the requirement for the memory bandwidth provided by the memory on the network adapter. On the other hand, some network

adapters such as Myrinet allow an integrated checksum and DMA capability, i.e., the DMA engine itself can compute the checksum of the data while performing the DMA operation. This avoids the checksum computation overhead on the host CPU and at the same time does not increase the memory bandwidth requirements on the network adapter (since it does not require any additional fetching of data).

*Exposing Scatter/Gather capabilities to the device driver:* Several network architectures including Myrinet, Quadrics and InfiniBand allow applications to send data from non-contiguous virtual address space and receive data into non-contiguous virtual address space, a capability often referred to as scatter/gather data transfer. For IP-based protocols this feature has several implications due to the functionality and implementation of these protocol stacks. For example, TCP/IP and UDP/IP copy the data from the application buffer into the kernel buffer and attach appropriate protocol headers for each segment of data before handing it over to the device driver. For such protocols, the header is generated inside the kernel. Implementing a zero-copy data transfer for such protocols would result in the data being present in the application buffer while the protocol header being generated in the kernel. To efficiently transfer these non-contiguous segments in a zero-copy manner by the kernel, a gather feature is quite important and required. Similarly, on the receiver side, the protocol headers need to be placed inside the kernel space while the data needs to be placed in the application buffer. In such a situation a scatter feature in the network adapter is quite beneficial.

Together with such a capability, we would also like to point out the requirement for scatter/gather capabilities with contiguous virtual address space but non-contiguous physical address space. Since, in our implementation, we retain the sockets API which allows arbitrary application buffers, the buffer from which the application tries to send data would be a contiguous virtual address space but does not have to be a contiguous physical address space. Therefore, the data transmission should be able to handle non-contiguous physical address space.<sup>2</sup> For handling such requirements in the data transmission, a scatter/gather capability for non-contiguous physical address space is very critical.

*Solicit bit for message segmentation:* The fine-grained pipelining described in this paper utilizes a feature known as solicit bit-based message segmentation. Using this feature, the Myrinet NIC allows several segments of data to be sent out separately but considered as a single message on the receiver side. The main advantage of this feature is that each of these segments does not need to have a separate protocol or MAC header; thus the available bandwidth can be utilized for transferring only the actual data. It is to be noted that since the Myrinet network is a cut-through net-

work, each segment does not need to have a MAC header to identify the destination node. As long as the solicit bit is set on the message segments, the Myrinet switch considers these segments to be a part of the same stream (message) and automatically forwards them to the same destination.

For store-and-forward networks, however, though this kind of an approach is possible, it might not be as beneficial for a large network. Store-and-forward networks require the entire link level message (MAC header + data specified by the MAC header) to be received by the switch before it can be forwarded to the destination node. Thus, there can be some amount of pipelining at the source nodes host and network adapter, but once the message reaches the first switch, there would be no difference between coarse-grained pipelining and fine-grained pipelining.

*Multiple DMA engine support:* Another important feature on the NIC that pipelining tries to exploit is the use of multiple DMA engines. The Myrinet NIC provides three DMA engines; one to transfer data from the host memory to the NIC memory and vice versa, one to transfer data from the NIC memory to the link and one to transfer data from the link to the NIC memory. It is to be noted that the current Myrinet network adapters are based on PCI and PCI-X I/O buses. These I/O buses are shared for data traffic in both directions (from host to NIC and from NIC to host). Thus, one DMA engine would be sufficient to perform the data transfer in both directions (SDMA and RDMA threads). On the other hand, the network link is full-duplex, i.e., data can be sent and received simultaneously. Thus, one dedicated DMA engine (SEND and RECV) is required for data transfer in each direction (from NIC memory to link and from link to NIC memory). It is to be noted that with the current PCI-Express technology network adapters will be capable of moving data from the host to the NIC and from the NIC to the host simultaneously in a dedicated manner. This can put a requirement for an additional DMA engine or full-duplex DMA capability to be present on the NIC.

*Memory bandwidth for internal data movement on the NIC:* As mentioned earlier, the Myrinet NIC performs checksum calculation during the DMA operation from the host memory to the NIC memory. Similarly, it performs a CRC computation during the DMA operation from the NIC memory to the link. This means that the only data accesses to the data in the NIC memory are by the DMA engines (one access while performing the DMA from the host memory to the NIC and one access while performing the DMA from the NIC memory to the link). The memory on the Myrinet NIC can support two accesses per clock cycle. On the other hand, the DMA engines can request one memory access each per clock cycle. So, two DMA engines can access the memory simultaneously at full speed. We utilized this to show fine-grained pipelining for uni-directional traffic. For bi-directional traffic on the other hand, the RECV DMA engine also needs to access the memory thus increasing the requirement for memory bandwidth. It is to be noted that for network architectures with full-duplex links and full-duplex

<sup>2</sup>It is to be noted that an approach, where the device driver issues multiple posts for the different physical pages is not possible since this would require the receiver to know the exact layout of the data on the sender side in order to post-appropriate-sized descriptors.

I/O buses (PCI-Express), the memory bandwidth requirements can be up to four times the link bandwidth in order to achieve non-blocked pipelining. Further, it is to be noted that this requirement is not specific to fine-grained pipelining, but a generic requirement to perform any kind of pipelining.

## 7. Related work

Several researchers have worked on implementing high-performance user-level sockets implementations over high-performance networks. Balaji et al. have worked on such pseudo-sockets layers over Gigabit Ethernet [7], GigaNet cLAN [8,39] and InfiniBand [6]. However, these sockets layers only deal with compatibility issues with existing applications and do not focus on compatibility with the existing IP infrastructure.

Similarly, there has been some previous work in the MPI domain. Madeleine MPI [5] and MPICH-G [28] are two implementations of MPI which focus on transparently handling heterogeneous network environments. However, these solutions are applicable for only MPI-based applications and are not usable by applications based on other common interfaces, such as sockets, DSM [3,29], Global Arrays [33], etc. On the other hand, since there exist sockets-based implementations for most of these environments, our implementation can be expected to provide a transparent compatibility with the applications as well as the IP-based infrastructure.

Trapeze [13] has implemented zero copy, checksum offloading, and a form of overhead pipelining based on TCP/IP on Myrinet. This research is notable in the sense that this was the one of the first to show that an IP-based protocol can achieve a significantly high performance. However, Trapeze provides a different API from the sockets interface and is not compatible with the traditional TCP/IP implementation. In our paper, we try to achieve a near theoretical performance while keeping the socket interface and compatibility existing the UDP/IP implementations.

For overhead pipelining, several studies have been done to achieve a middle-grained pipelining [24,36,45], which splits a large packet into smaller sized segments so that the overheads for a packet at a layer and the same packet at the next layer are partly pipelined. The middle-grained pipelining can achieve a lower latency than coarse-grained pipelining but is not as efficient as fine-grained pipelining. Further, since middle-grained pipelining splits a packet into separate segments, each with a separate segment header, it sacrifices some of the bandwidth.

A notable research on formal verification of the NIC firmware is performed by Kumar et al. [31], which employs a model checking approach to implement the NIC firmware. They use Event-driven State-machine Programming (ESP); a language for writing firmware for programmable devices to verify the retransmission protocol, memory safety, and the deadlock free property of the VMMC [15] firmware.

We expect that their approach can effectively help implement our suggested model.

## 8. Concluding remarks and future work

While a number of user-level protocols have been developed to reduce the gap between the performance capabilities of the physical network and the performance actually available, their compatibility issues with the existing sockets-based applications and IP-based infrastructure has been an area of major concern. To address these compatibility issues while maintaining a high performance, a number of researchers have been looking at alternative approaches to optimize the existing traditional protocol stacks. Broadly, previous research has broken up the overheads in the traditional protocol stack into four related aspects, namely: (i) compute requirements and contention, (ii) memory contention, (iii) I/O bus contention and (iv) system resources' idle time.

There has been some previous research which deals with some of these aspects. For example, POEs have been recently proposed as an industry standard for offloading the compute intensive components in protocol processing to specialized hardware. However, these approaches require network adapters supported with specialized ASIC-based chips which implement the protocol processing and are not generic enough to be implemented on most network adapters. Further, these deal with only the compute requirement and memory contention issues and do not address the remaining issues. In short, to the best of our knowledge, there is no work which deals with all these issues in an integrated manner while maintaining backward compatibility with existing applications and infrastructure.

In this paper, we address each of these issues and propose solutions for minimizing these overheads. We also modify the existing UDP/IP implementation over Myrinet to demonstrate the capabilities of these solutions. We first utilize the earlier proposed techniques to implement a partial offload of the UDP/IP protocol stack to address the first two issues, i.e., compute requirements and memory contention. Next, we modify the device driver to allow a delayed posting of descriptors in order to reduce the contention at the I/O bus between descriptor posting and the DMA operations of the actual outgoing or incoming data. Finally, we implement a fine-grained pipelining technique on the firmware of the network adapter to minimize the link idle time in order to achieve a high performance. Further, all these enhancements to the UDP stack are completely compatible not only with existing applications and infrastructure, but also with the existing UDP implementations. Our experimental results show that with our implementation of UDP, termed as E-UDP, can achieve up to 94% of the theoretical maximum bandwidth. We also present a mathematical performance model which allows us the study the performance of our design for various system architectures and network speeds.

Reliability is a critical feature for running several real applications over a cluster-of-clusters environment. There has



been some previous work related to reliability over UDP [10,18]. Such reliability relieves us of the requirement of a heavy protocol such as TCP/IP and allows us to achieve high performance. As a part of our future work, we plan to integrate this kind of reliability into E-UDP. This would allow us to deploy our solution in a real cluster-of-clusters environment and measure the actual application level performance gains achievable through this solution. Also, we are currently looking into offloading the per-packet overheads in the current UDP/IP stack on to the network adapter. This would not only allow us to achieve a better performance for small messages (where the per-packet overhead is dominant), but also a better scalability for faster networks.

## Acknowledgments

This research is supported in part by Department of Energy's Grant #DE-FC02-01ER25506, National Science Foundation's Grants #CCR-0204429, and #CCR-0311542 and the Post-doctoral Fellowship Program of Korea Science and Engineering Foundation (KOSEF). The authors thank Sundeep Narravula and Karthikeyan Vaidyanathan for helping to improve the performance model.

## References

- [1] Ammasso, Inc., <http://www.ammasso.com>
- [2] Chelsio Communications, <http://www.chelsio.com>
- [3] C. Amza, A. Cox, S. Dwarkadas, P. Keleher, H. Lu, R. Rajamony, W. Yu, W. Zwaenepoel, Treadmarks: shared memory computing on networks of workstations, *IEEE Comput.* 29 (2) (1996) 18–28.
- [4] InfiniBand Trade Association. <http://www.infinibandta.org>
- [5] O. Aumage, G. Mercier, MPICH/MADIII: a cluster of clusters enabled MPI implementation, in: *The Proceedings of Third International Symposium on Cluster Computing and the Grid*, May 2003.
- [6] P. Balaji, S. Narravula, K. Vaidyanathan, S. Krishnamoorthy, J. Wu, D.K. Panda, Sockets direct protocol over infiniband in clusters: is it beneficial?, in: *The Proceedings of ISPASS 2004*, Austin, Texas, March 10–12, 2004.
- [7] P. Balaji, P. Shivam, P. Wyckoff, D.K. Panda. High performance user level sockets over Gigabit Ethernet, in: *The Proceedings of Cluster'02*, Chicago, Illinois, September 23–26, 2002, pp. 179–186.
- [8] P. Balaji, J. Wu, T. Kurc, U. Catalyurek, D.K. Panda, J. Saltz. Impact of high performance sockets on data intensive applications, in: *The Proceedings of HPDC-12*, Seattle, Washington, June 22–24, 2003, pp. 24–33.
- [9] N.J. Boden, D. Cohen, R.E. Felderman, A.E. Kulawik, C.L. Seitz, J.N. Seizovic, W.K. Su, Myrinet: a gigabit-per-second local area network, *IEEE-Micro* 15 (1) (February 1995) 29–36.
- [10] T. Bova, T. Krivoruchka, Reliable UDP protocol, in: *Internet Draft*, February 1999.
- [11] P. Buonadonna, A. Geweke, D.E. Culler, BVIA: an implementation and analysis of virtual interface architecture, in: *Proceedings of Supercomputing*, 1998.
- [12] J. Chase, A. Gallatin, A. Lebek, Y.G. Yocum, Trapeze API, Technical Report CS-1997-21, Duke University, Durham, NC, November 1997.
- [13] J. Chase, A. Gallatin, K. Yocum, End-system optimizations for high-speed TCP, *IEEE Commun. Mag.* 39 (4) (2001) 68–75.
- [14] H.J. Chu, Zero-copy TCP in solaris, in: *Proceedings of 1996 Winter USENIX*, 1996.
- [15] C. Dubnicki, A. Bilas, Y. Chen, S. Damianakis, K. Li, VMMC-2: efficient support for reliable, connection oriented communication, in: *Proceedings of Hot Interconnects*, 1997.
- [16] W. Feng, J. Hurwitz, H. Newman, S. Ravot, L. Cottrell, O. Martin, F. Coccetti, C. Jin, D. Wei, S. Low, Optimizing 10-gigabit Ethernet for networks of workstations, clusters and grids: a case study, in: *Proceedings of ICS '03*, Phoenix, Arizona, November 2003.
- [17] Message Passing Interface Forum, MPI: a message-passing interface standard, <http://www.mpi-forum.org>, May 1994.
- [18] E. He, J. Leigh, O. Yu, T. DeFanti, Reliable blast UDP: predictable high performance bulk data transfer, in: *the Proceedings of the IEEE International Conference on Cluster Computing*, 2002.
- [19] G.J. Holzmann, *Design and Validation of Computer Protocols*, Prentice-Hall, Englewood Cliffs, NJ, 1991.
- [20] G.J. Holzmann, The Model Checker SPIN, *IEEE Trans. Software Eng.* 23 (5) (1997) 279–295.
- [21] J. Hurwitz, W. Feng, End-to-end performance of 10-gigabit Ethernet on commodity systems, *IEEE Micro.* 24 (1) (2004) 10–22.
- [22] IEEE. IEEE Std. 802.3ae-2002, Media Access Control (MAC) Parameters, Physical Layers, and Management Parameters for 10 Gb/s Operation, August 2002.
- [23] Myricom Inc., The GM Message Passing System, <http://www.myri.com>, January 2000.
- [24] H.A. Jamrozik, M.J. Feeley, G.M. Voelker, J.E. II, A.R. Karlin, H.M. Levy, M.K. Vernon, Reducing network latency using subpages in a global memory environment, in: *Proceedings of ASPLOS-VII*, October 1996.
- [25] H.-W. Jin, P. Balaji, C. Yoo, J.Y. Choi, D.K. Panda, Exploiting NIC architectural support for enhancing IP based protocols on high performance networks, Technical Report OSU-CISRC-5/04-TR37, The Ohio State University, Columbus, OH, May 2004.
- [26] H.-W. Jin, K.S. Bang, C. Yoo, J.Y. Choi, H.J. Cha, Bottleneck analysis of a gigabit network interface card: formal verification approach, in: *Proceedings of the Ninth International SPIN Workshop on Model Checking of Software (SPIN 2002)*, Lecture Notes in Computer Science, vol. 2318, Springer, Berlin, April 2002.
- [27] H.-W. Jin, C. Yoo, Latency analysis of UDP and BPI on Myrinet, in: *Proceedings of the 18th IEEE International Performance, Computing, and Communication Conference (IPCCC'99)*, February 1999, pp. 185–191.
- [28] N. Karonis, B. Toonen, I. Foster, MPICH-G2: a grid-enabled implementation of the message passing interface, *J. Parallel Distributed Comput.* 63 (5) (2003).
- [29] P. Keleher, A.L. Cox, S. Dwarkadas, W. Zwaenepoel, Distributed shared memory on standard workstations and operating systems, in: *The Proceedings of the 1994 Winter USENIX Conference*, January 1994.
- [30] J.S. Kim, K. Kim, S.I. Jung, SOVIA: a user-level sockets layer over virtual interface architecture, in: *The Proceedings of Cluster '01*, California, USA, October 8–11, 2001, pp. 399–408.
- [31] S. Kumar, K. Li, Using model checking to debug device firmware, in: *Proceedings of OSDI*, 2002.
- [32] Z. Manna, A. Pnueli, *The Temporal Logic of Reactive and Concurrent Systems*, Springer, Berlin, 1992.
- [33] J. Nieplocha, R.J. Harrison, R.L. Littlefield, Global Arrays: a Portable "Shared Memory" programming model for distributed memory computers, in: *The Proceedings of Supercomputing*, 1994.
- [34] S. Pakin, M. Lauria, A. Chien, High performance messaging on workstations: Illinois Fast Messages (FM), in: *Proceedings of Supercomputing*, 1995.
- [35] F. Petrini, W. C. Feng, A. Hoisie, S. Coll, E. Frachtenberg, The quadrics network (QsNet): high-performance clustering technology, in: *The Proceedings of Hot Interconnects 9*, August 2001.
- [36] L. Prylli, R. Westerlin, B. Tourancheau, Modeling of a high speed network to maximize throughput performance: the experience of BIP over Myrinet, in: *Proceedings of PDPTA '98*, 1998.

- [37] G. Regnier, D. Minturn, G. McAlpine, V.A. Saletore, A. Foong, ETA: experience with an Intel Xeon Processor as a packet processing engine, *IEEE Micro* 24 (1) (2004) 24–31.
- [38] H.V. Shah, D.B. Minturn, A. Foong, G.L. McAlpine, R.S. Madukkarumukumana, G.J. Regnier, CSP: a novel system architecture for scalable internet and communication services, in: *The Proceedings of the Third USENIX Symposium on Internet Technologies and Systems*, San Francisco, CA, March 2001, pp. 61–72.
- [39] H.V. Shah, C. Pu, R.S. Madukkarumukumana, High performance sockets and RPC over virtual interface (VI) architecture, in: *The Proceedings of the CANPC workshop (held in conjunction with HPCA)*, 1999, pp. 91–107.
- [40] P. Shivam, P. Wyckoff, D.K. Panda, EMP: Zero-copy OS-bypass NIC-driven gigabit Ethernet message passing, in: *The Proceedings of ICS*, Denver, Colorado, November 10–16, 2001, pp. 57–64.
- [41] P. Shivam, P. Wyckoff, D.K. Panda, Can User-Level protocols take advantage of multi-CPU NICs?, in: *The Proceedings of IPDPS*, Fort Lauderdale, Florida, April 15–19, 2002.
- [42] Y. Turner, T. Brecht, G. Regnier, V. Saletore, G. Janakiraman, B. Lynn, Scalable networking for next-generation computing platforms, in: *Proceedings of SAN*, 2004.
- [43] USNA. TTCP: a test of TCP and UDP performance, December 1984.
- [44] T. von Eicken, A. Basu, V. Buch, W. Vogels, U-Net: a user-level network interface for parallel and distributed computing, in: *Proceedings of SOSP*, December 1995.
- [45] R.Y. Wang, A. Krishnamurthy, R.P. Martin, T.E. Anderson, D.E. Culler, Modeling and optimizing computation pipelines, in: *Proceedings of SIGMETRICS*, June 1998.
- [46] E. Yeh, H. Chao, V. Mannem, J. Gervais, B. Booth, Introduction to TCP/IP Offload Engine (TOE), <http://www.10gea.org>, May 2002.
- [47] C. Yoo, H.-W. Jin, S.C. Kwon, Asynchronous UDP, *IEICE Trans. Commun.* E84-B (12) (December 2001) 3243–3251.

**Hyun-Wook Jin** is a Research Associate of Computer Science and Engineering at The Ohio State University. His research interests include cluster and grid computing, operating systems, high-speed network protocols, and network interface card architectures. He has received the Ph.D., M.S., and B.S. degrees in computer science and engineering from Korea University, Seoul, Korea, in 2003, 1999, and 1997. He is a member of IEEE.

**Pavan Balaji** is a Ph.D. student in the Computer Science and Engineering Department at The Ohio State University. His research interests include high-speed interconnects, efficient IP-based protocols (TCP/IP, UDP/IP, iWARP) and the sockets programming interface for cluster and grid computing. He has more than 20 publications in these areas. He has a B.Tech. in Computer Science and Engineering from The Indian Institute of Technology at Madras, India. He is a member of IEEE.

**Chuck Yoo** received the B.S. degree in electronics engineering from Seoul National University, Seoul, Korea and the M.S. and Ph.D. in computer science in University of Michigan. He worked as a researcher in Sun Microsystems Lab. from 1990 to 1995. He joined the Computer Science and Engineering Department, Korea University, Seoul, Korea in 1995, where he is currently a Professor. His research interests include high performance network, multimedia streaming, and operating systems.

**Jin-Young Choi** received the B.S. degree from Seoul National University, Seoul, Korea, in 1982, the M.S. degree from Drexel University in 1986, and the Ph.D. degree from University of Pennsylvania, in 1993. He is currently a Professor of Computer Science and Engineering Department, Korea University, Seoul, Korea. His current research interests are in real-time computing, formal methods, programming languages, process algebras, security, software engineering, and protocol engineering.

**Dhableswar K. Panda** is a Professor of Computer Science and Engineering at The Ohio State University. His research interests include parallel computer architecture, high performance networking, and network-based computing. He has published over 170 papers in these areas. His research group is currently collaborating with National Laboratories and leading companies on designing various communication and I/O subsystems of next generation HPC systems and datacenters with modern interconnects. The MVAPICH (MPI over VAPI for InfiniBand) package developed by his research group (<http://nowlab.cis.ohio-state.edu/projects/multi-iba/>) is being used by more than 230 organizations world-wide to extract the potential of InfiniBand-based clusters for HPC applications. Dr. Panda is a recipient of the NSF CAREER Award, OSU Lumley Research Award (1997 and 2001), and an Ameritech Faculty Fellow Award. He is a senior member of IEEE Computer Society and a member of ACM.