

On Network CoProcessors for Scalable, Predictable Media Services *

Raj Krishnamurthy, Karsten Schwan

Richard West

Center for Experimental Research in Computer Systems

Department of Computer Science

Georgia Institute of Technology

Boston University

Atlanta, GA 30332

Boston, MA

{rk, schwan}@cc.gatech.edu

richwest@cs.bu.edu

Marcel-Cătălin Roșu

IBM T.J. Watson Research Center

mrosu@watson.ibm.com

October 24, 2002

*This work was supported in part by the Department of Energy under its NGI program, by the National Science Foundation under a grant from the Division of Advanced Networking Infrastructure and Research, and by hardware/software donations from Intel Corporation and WindRiver Systems

Abstract

This paper presents the embedded realization and experimental evaluation of a media stream scheduler on Network Interface (NI) CoProcessor boards. When using media frames as scheduling units, the scheduler is able to operate in real-time on streams traversing the CoProcessor, resulting in its ability to stream video to remote clients at real-time rates. The contributions of this paper are its detailed evaluation of the effects of placing application- or kernel-level functionality like packet scheduling on NIs rather than the host machines to which they are attached. The main benefits of such placement are (1) that traffic is eliminated from the host bus & memory subsystem, thereby allowing increased host CPU utilization for other tasks, and (2) that NI-based scheduling is immune to host-CPU loading, unlike host-based media schedulers that are easily affected even by transient load conditions. An outcome of this work is a proposed cluster architecture for building scalable media servers, by distributing schedulers and media stream producers across the multiple NIs used by a single server and by clustering a number of such servers using commodity network hardware and software.

KEYWORDS: cluster machines, multimedia services, embedded systems, quality of service, operating systems, real-time systems, data streaming, packet scheduling

1 Introduction

Background. The scalable delivery of media and web services to end users is a well-recognized problem. At the network level, researchers have developed multicast techniques[4], media caching or proxy servers[2], reservation-based communication services[16] and specialized media transmission protocols [25]. For server hardware, scalability is sought by using extensible SMP and cluster machines[5, 11]. Scalability for server software is attained by using dynamic load balancing across parallel/distributed server resources[50], and by using admission control and online request scheduling[51, 66] for CPUs[41, 9, 24], network links[63, 62, 66], and disks[11]. Complementary to such work are application-level or end-to-end solutions[35, 46] that adapt server and/or client behavior in response to changes in users' Quality of Service (QoS) needs and in resource availability [59, 47].

Scalable Cluster Services. This paper addresses the scalability of media servers. Its approach utilizes servers constructed as clusters of processor/storage nodes, each of which has a network interface processor (NI) linking it to the cluster's high performance system area network[60]. A unique aspect of the approach is that each NI has a

programmable CoProcessor, with local memory, direct access to devices, and the ability to run both basic protocol processing and certain application-specific tasks. The CoProcessors explored to date include ATM FORE[49], Myrinet[60], I2O-compliant network interface boards[22, 38, 21], and most recently, gigabit ethernet attached to IXP1200[55, 52, 18, 68] boards. We have also explored Xilinx FPGA Coprocessor boards[65] in servers with gigabit links and the ability of the FPGA Coprocessor to meet the packet-time requirements of 10Gbps links [30]. The server hardware configuration used in this paper is from a generation of CoProcessors, (developed for the I2O standard), equipped with the i960 RD Coprocessor (our systems software does not use any I2O standard-style driver partitioning or any of the hardware registers in the I2O hardware units). This research simply uses the i960 RD processor on the NI with associated PCI bridge chipsets, as a i960 RD Coprocessor. These NIs have two 100Mbps Ethernet links, a PCI interface to the host CPU, and two SCSI interfaces directly attached to disk devices. These CoProcessors are attached to a prototypical server system comprised of 16 quad Pentium Pro nodes, where host-to-host communications are supported by VxWorks TCP/IP stack protocols (like TCP and UDP), and where media streams may flow from server disks to clients via host CPUs or directly via the i960RD boards [38, 21].

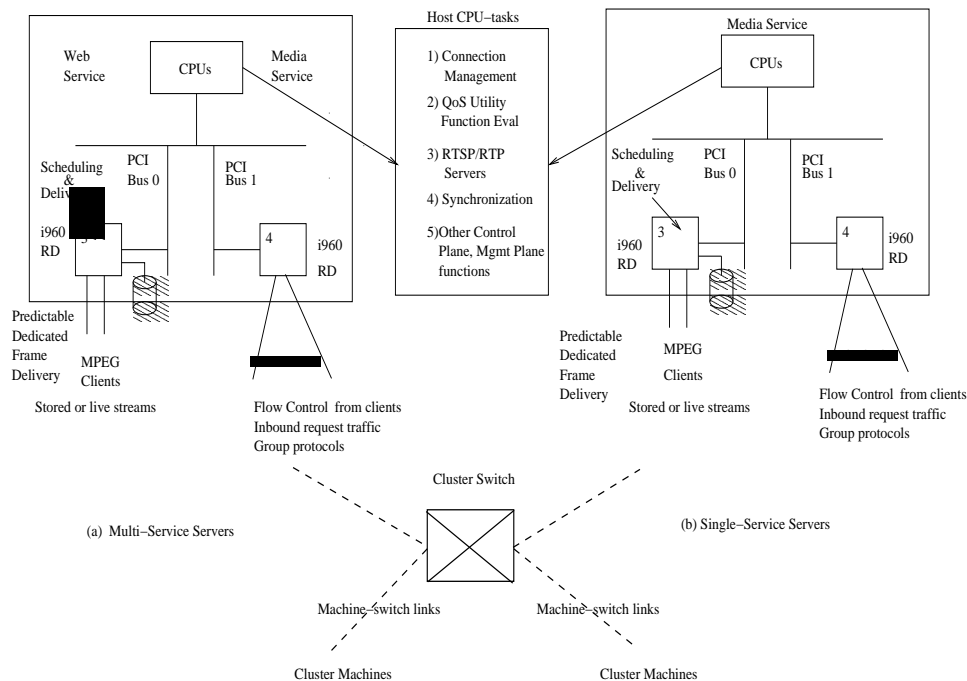


Figure 1: Cluster Hardware: Host CPU, NI CoProcessor and Interconnect

The approach to server organization we advocate is one that views a server like the one depicted in Figure 1 as

an information processing, storage, and delivery engine that is programmed at two levels of abstraction, reflecting the hardware configuration being employed:

1. Low-level services execute on the CoProcessors (NIs), as demonstrated in Figure 1 for the concrete examples of frame scheduling and predictable frame delivery for streaming media services. Whether executed synchronously or asynchronously with application programs' communications, these services may be viewed as application-specific *extensions* supported by a *runtime* resident on the NI.
2. Higher level services run on host nodes but they may utilize CoProcessor services like media scheduling via an explicit NI-provided interface. Figure 1 shows sample media service tasks that are run on host CPUs, with predictable frame delivery dedicated to run on NIs.

Contributions. Our approach to improving the scalability of media servers is to amplify their capabilities by extension of their communication CoProcessors with application-specific functionality. The intent is to use cluster nodes with their complete OS functionality, large memories, and deep cache hierarchies for the computational and data management tasks for which they are well-suited, while using NIs for communication-centric tasks like packet scheduling. The following insights are gained by this experimental research.

1. Efficient execution on standard NIs. Even NIs like Intel's i960RD-based boards [38, 21] are feasible platforms on which to run application-specific extensions of communication functionality. Specifically, performance-critical communication extensions can be executed with high performance on such COTS (Commercial Off-The-Shelf) CoProcessors running with standard operating system software (e.g., VxWorks[64]). This is encouraging given the fast-moving nature of NI hardware development and the onerous task of porting and re-reporting the specialized runtime software used in previous work [45, 17].

2. Improved predictability for NI-based streaming services. On host CPUs, the time-critical execution of services is easily jeopardized by the need to run a mix of applications. Specifically, streaming services like media schedulers running on host-CPU are easily affected even by transient loading conditions, whereas media schedulers running directly on NIs are immune to such host-CPU loading. The key result is that CoProcessor-based scheduling substantially improves the jitter experienced by streaming services. Such improvements are particularly important for real-time media services like remote surveillance or telepresence[56, 33].

3. Improved server scalability by use of extended NIs. Services that frequently execute *device-* or *network-near* functions are known to run faster on CoProcessors, because their execution does not involve I/O busses, host

memory, and host CPUs. This paper presents results that better quantify the load reductions – termed *traffic elimination* – on such node resources derived from NI-based service execution, thus freeing up these resources for use by other server tasks. Traffic elimination is attained by implementing stream-selective lossiness in overload conditions via window- and time-constrained scheduling of MPEG video frames, employing the DWCS (Dynamic Window Constrained Scheduling) algorithm [63, 62]. Packet scheduling also serves to guarantee differential packet rates and deadlines to meet clients’ individual QoS needs, and to eliminate traffic.

Summary of Experimental Results. A DWCS-based media scheduler can run almost as fast on a relatively slow CoProcessor as its corresponding implementation on a standard workstation platform. Specifically, the scheduling latency of the host-based implementation of DWCS reported in [63, 62] is $\approx 50\mu s$ on an Ultra Sparc CPU (300 MHz) with quiescent load. In comparison, the scheduler’s execution on a 66 MHz i960 RD CoProcessor is $\approx 67\mu s$, despite the fact that the i960 RD is roughly 4 times slower than the Ultra Sparc host CPU. Reasons for high scheduler performance on the CoProcessor are presented in Section 3.2.

Performing packet scheduling on the NI rather than the host reduces load on the server’s I/O busses, CPU, and memory resources. In one experiment described in Section 3.3, media content are sent directly from the NI to a remote host, thus eliminating $\approx 132\text{MB/sec}$ of data load from PCI bus and memory of the host node for a 32-bit wide bus running at 33MHz. With multiple media streams, the load eliminated consumes a substantial fraction of the machine’s available PCI bandwidth. The importance of reducing PCI load is evident when considering that a media stream that fully utilizes 1Gbps links would consume 25% of the capacity of a 64bit/66MHz PCI interface, which shows that multi-gigabit interconnects (e.g., 10Gbps ethernet) would consume more than the capacity of the fastest currently available 64bit/66MHz PCI interfaces.

Section 3.4 demonstrates the comparatively high predictability of CoProcessor- vs. host-based media services, by showing that transient loads can substantially affect host-CPU based schedulers, observing performance degradation for media scheduling even for relatively low CPU utilizations (i.e., 45%) and severe degradation for high CPU utilizations (i.e., over 60%).

Previous Work. Earlier work conducted by our group did not address scalable cluster media services. Further, rather than using standard host-CoProcessor interfaces, we developed and experimented with zero-copy interfaces and with a software architecture for realizing a rich set of communication instructions on NIs, termed a *(Distributed) Virtual Communication Machine* (DVCM[45]). The idea was to permit application programs to dynamically extend the current set of communication instructions resident on the NI to support their specific

needs[45, 48], as subsequently also done in the SPINE project[17]. As with SPINE and other prior work on dynamically extending operating system kernels[14, 8], the services implemented by the DVCM can vary over time, in keeping with the needs of current cluster applications. The prototype described in this paper does not re-implement all of DVCM. Instead, our goals are (1) to gain an understanding of how to build a scalable media server from clustered node-NI pairs, (2) to evaluate in detail an NI's ability to support multimedia scheduling services, and (3) to compare NI- vs. host-based scheduling service realizations. Furthermore, while an earlier publication by our group already demonstrated the basic ability of a CoProcessor to perform media scheduling[29], this paper provides basic insights into the benefits of using CoProcessors or hosts, including an analysis of the overheads and tradeoffs concerning media scheduling with respect to locating media schedulers on CoProcessors vs. hosts, using multiple vs. single CoProcessors. In addition, we evaluate the utility of certain CoProcessor hardware, such as the benefits derived from caches, software floating-point and specialized scheduler hardware units.

2 Network CoProcessor-Based Media Scheduling

2.1 Software Architecture of NI-based Application Services

Our NI-based support for application-specific services is structured as three sets of software modules: host interface, runtime support, and application-specific extensions.

Host Interface. The interface functions exported by the NI to the host make NI-resident communication services appear to application programs as specialized 'communication instructions'. These instructions are accessible via memory-mapped pages shared by a host-resident application and the NI-resident media scheduling service. Pages contain control information as well as the communication buffers used for message transfers from NI to host and vice versa, much like the efficient message interfaces used in high performance messaging systems like FM [40]. With this interface, media frame producers, running as application threads, may stream frames to remote clients using the NI-resident scheduler. Remote consumers may forward frames to other consumers or buffer frames for display or storage, where frames are scheduled for delivery by the NI-resident scheduler.

Runtime. This set of NI-based modules supports the implementation of application-specific NI functionality. Using the VxWorks real-time operating system[64] as a basis, additional functionality to exploit this NI's specific hardware and to implement efficient media scheduling include a fixed-point library for efficient implementation of certain scheduling computations, driver front-ends to initialize controllers/storage, timestamp counter rollover

management, and circular queues and heaps as the buffer structures used for media frames.

Extensions. NI extensions support specific applications’ needs. Examples evaluated in our previous work include atomic read and write operations for remote NIs and efficient implementations of cluster-wide synchronization operations[45, 48]. This paper implements and evaluates the DWCS scheduler for streaming media on the i960 RD cards.

2.2 A CoProcessor-based Media Scheduling Service

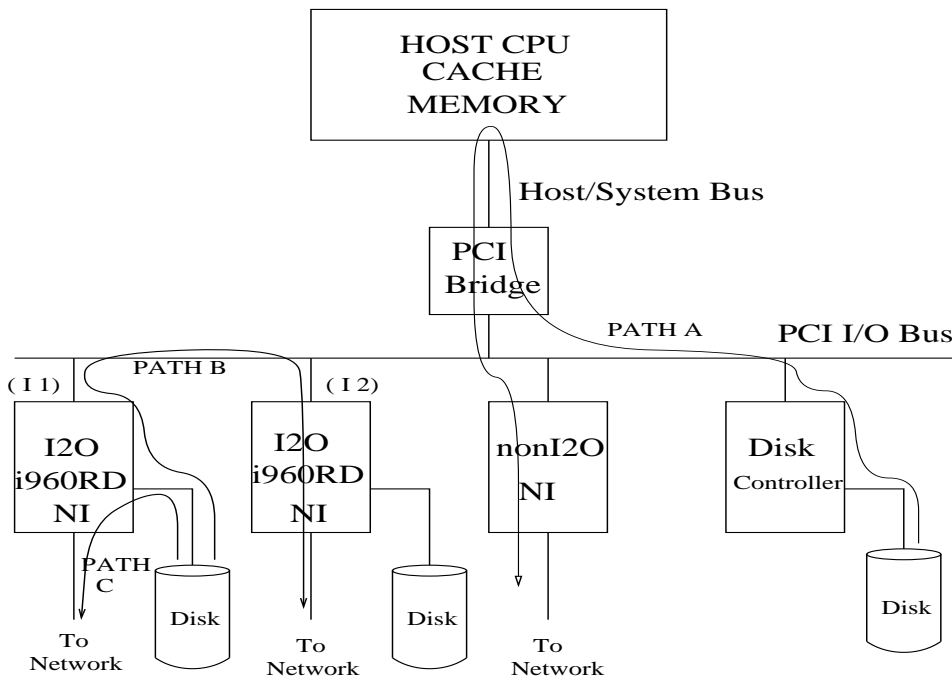


Figure 2: Frame Transfer Paths.

Alternative Service Configurations and Basic Performance Factors. The performance of a media scheduling service is strongly affected by the distribution of streams and schedulers across the underlying cluster node/NI pairs. Alternative stream and scheduler configurations are depicted in Figure 2, where multiple NIs (in this case i960 RD cards) are present on each cluster node’s I/O bus. One or more of these NIs may run the media scheduler and also support disks directly attached to them. This results in three possible paths traversed by the media frames being scheduled, shown in Figure 2 as Paths A, B and C.

Path A represents the case where media scheduling is performed only on host CPUs, which also implies that all

media streams must touch upon hosts. As a result, frames sent to clients by the server are transferred from a server disk attached to a SCSI controller card, to server host memory (via a PCI interconnect), then again transferred via the PCI interconnect to a non-i960RD NI and finally, sent via the network to whichever clients requested the media stream. Stated in more detail, the server's node OS transfers the MPEG file from disk to its filesystem buffer cache, then to the application-level thread that has opened the MPEG file for reading. This involves at least two memory copies as well as the traversal of memory hierarchies and of bus domains (i.e., I/O bus to system bus). After having been scheduled, the second part of Path A involves the transfer of frames from host CPU memory to the network via the NI, again involving multiple memory hierarchy and bus traversals.

Path B depicts a configuration in which a media stream originates either on another machine or on disks directly attached to an i960 RD card (I1), and is then sent to some client via a second i960 RD card (I2) running the media scheduler. This path represents the general case of NI-based media stream scheduling, involving multiple NIs attached to one server node, with each NI specialized to perform certain tasks. This path involves the I/O bus, but completely eliminates the uses of host CPU or memory.

Finally, the 'best' case in terms of host node resource usage is depicted in Path C, where a single NI acts as both the source of the media stream using a disk attached to it and also executes the media scheduler. Compared to Path A, this path eliminates uses of the I/O bus, host CPU and memory. While imposing additional load onto the NI, one potential advantage of this configuration is the relative 'closeness' of the media scheduler to the network, which facilitates rapid configuration in response to changes in network behavior. A sample dynamic reconfiguration is one that adjusts the scheduler's degree of lossiness in response to observing changes in the number of packet retransmissions currently experienced for this media stream. Another possible advantage of this configuration is the ability to share a single set of buffers between the NI's disk and network interface, thereby avoiding additional memory copies[39].

Summarizing the performance characteristics of different service configurations:

1. Performance is impacted both by the speed of scheduling actions and by the total percentage of a media stream's data that traverses a cluster node's memory hierarchies.
2. It is affected by the total protocol processing overheads experienced on NI and/or nodes, and by the disk access overheads experienced for media streams originating at or destined to storage devices.
3. It also depends on the number of bus-domain traversals (system bus to PCI bus and vice-versa) experienced

by streams.

4. Scheduling and streaming performance are also affected by the presence of other programs executing on nodes and NIs, measurable not only as changes in throughput, but also in terms of the delay-jitter and loss experienced for media streams and scheduling actions.

Host-based Scheduler Implementation. Our media scheduler is based on the DWCS (Dynamic Window-Constrained Scheduling) algorithm described and evaluated in [63, 62]. While shown to be a highly efficient algorithm, most important to this paper are the ways in which the performance of this scheduling service is affected by alternative implementation choices for its host- vs. NI-based realizations. Factors governing performance include the respective hardware and operating system platforms, certain hardware features, and memory structures. This large variety of factors precludes analytical comparisons of scheduling performance based on the alternative paths shown in Figure 2. For instance, the media scheduler's implementation on the host CPU running SUN's Solaris operating system is embedded into a separate scheduling process. By presenting a shared memory-based API (using System V shared memory) to processes that generate media content, this media scheduler is made independent of the source of media data (e.g., from host-attached disks or from remote nodes), many media streams can be scheduled by a single host-resident scheduler, and packets in any stream may be enqueued in scheduler data structures, concurrent with scheduling analysis and dispatch of previously enqueued packets. Additional scalability in scheduler operation with respect to number of packet streams, stream rates and 'tightness' of deadlines and loss-tolerance is achieved by varying the rates of scheduler actions[63, 62].

The NI-based Scheduler Implementation. The comparatively lighter-weight NI-embedded implementation of the DWCS scheduler is shown in Figure 3. It is layered on top of the NI's VxWorks real-time operating system. It runs as a single VxWorks thread, uses pinned memory for disk and/or network buffers and to interact with other NI-resident threads. Its host interface maps some of its memory to the host via the PCI device. Compact data structures (scheduler attributes or scheduler frame descriptors) for packet schedule representation minimize the use of NI memory, and memory usage is reduced further by not copying frames, whenever possible. These implementation choices attempt to compensate for the relative resource paucity on the NI. In addition, the NI-based scheduler takes advantage of certain hardware features existing on the i960 RD card. First, the code used for scheduling analysis is decoupled from the schedule representations (ie., scheduling data structures). The intent is to evaluate the performance effects of using alternative representations, including those that use the i960RD's hardware-supported FCFS circular buffer queues. This is important because with DWCS scheduling, packets

in a given stream (at the same priority level) may be scheduled based on a service tag associated with each packet. Second, by using one thread for both packet scheduling and dispatch, a single data structure can hold all frame descriptors (or the attributes describing them), thus conserving memory. Also, packets will not experience additional queuing delay and jitter in dispatch queues[63, 62].

The NI-based scheduler operates as follows. It is booted in conjunction with the VxWorks Operating System, from flash-ROM on the i960 RD NI card. Initialization code in the kernel spawns the scheduler thread. Any media frames received by the NI-based scheduler are temporarily stored on the NI, using the i960 RD NI's 4MB of on-board memory, which may be expanded to 36MB. To conserve memory, only a single copy of 'to-be-scheduled' frames is kept in NI memory, and scheduling analysis and dispatch directly manipulate addresses of frames. Frames are stored on a per-stream basis. Head-of-line packets in each stream form loss-tolerance and deadline heaps and encode stream priority values. The scheduler must pick the stream with the highest priority according to rules described in [63, 62].

Storing frames directly in NI memory (rather than in host memory) reduces the overall scheduling analysis and dispatch latency for each frame and the jitter experienced by a sequence of frames. It also reduces mean frame queuing delay, since frames need not be 'pulled' from host memory. A detailed analysis of the performance effects of locating various data structures on NIs vs. hosts appears elsewhere [45]. As shown in Figure 3, a

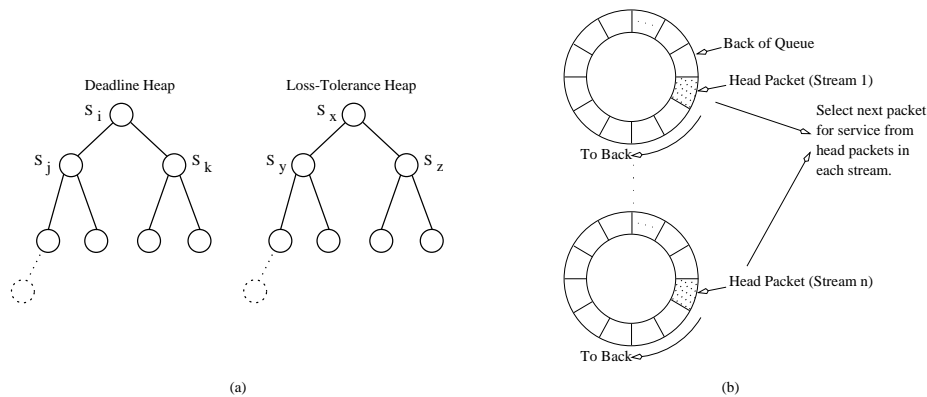


Figure 3: (a) The NI implementation of DWCS uses two heaps: one for deadlines and another for loss-tolerances. (b) Using a circular queue for each stream eliminates the need for synchronization between the scheduler that selects the next packet for service, and the server that queues packets to be scheduled.

circular buffer is maintained for each stream, with separate head and tail pointers. Frame producers inject frames into the scheduler using the tail pointer, and the scheduler reads frames using the head pointer. This eliminates any

explicit synchronization needs between readers and writers. A deeper understanding of the scheduler’s operation requires knowledge of the algorithm it executes, which is outlined next.

Some Details about the DWCS Scheduling Algorithm. DWCS is designed to maximize network bandwidth usage in the presence of multiple packets that have individual delay constraints and loss-tolerances. The per-packet delay and loss-tolerance attributes are derived from higher-level application constraints:

- *Deadline* – this is the latest time a packet can *commence* service. It is determined from a specification of the maximum allowable time between servicing consecutive packets in the same stream.
- *Loss-tolerance* – this is specified as a value x_i/y_i , where x_i is the number of packets that can be lost or transmitted late for every *window*, y_i , of consecutive packet arrivals in the same stream, i . For every y_i packet arrivals in stream i , a minimum of $y_i - x_i$ packets must be scheduled on time, while at most x_i packets can miss their deadlines and be either dropped or transmitted late, depending on whether or not the attribute-based QoS for the stream allows some packets to be lost.

At any time, all packets in the same stream have the same loss-tolerance, while successive packets in the same stream have deadlines that are offset by fixed amounts from their predecessors.

Using these attributes, DWCS has the following abilities: (1) it can limit the number of late packets over finite numbers of consecutive packets in loss-tolerant or delay-constrained, heterogeneous traffic streams; (2) it does not require a-priori knowledge of the worst-case loading from multiple streams to establish the bandwidth allocations necessary to meet per-stream delay and loss-constraints; (3) it can safely drop late packets in lossy streams without unnecessarily transmitting them, thereby avoiding needless bandwidth consumption; and (4) it exhibits both fairness and unfairness properties when necessary. Proofs of these properties and additional detail about DWCS appear in [63, 62]. [61] describes an approach where deadlines are compared before window-constraints. This approach was intended to provide hard guarantees on loss constraints. The slight variations of the DWCS algorithm in [61] currently being investigated do not concern its running time or performance and are therefore, not relevant to this paper’s contents, which is based on the original version of DWCS [63, 62].

Discussion. One property of DWCS (and of other media schedulers) is the fact that packets are dropped under certain conditions. This means that a scheduler operating on the NI will typically not forward all packets to the host and/or to other recipients. The resulting reduction in node resource usage is one motivation for placing scheduling onto the NI. Another motivation is the network-near nature of an NI-resident media scheduler, which enables it to rapidly change its operation in response to changes in network conditions. In previous work, we

have demonstrated performance advantages due to such network nearness for other NI-resident communication services[48]. Furthermore, in previous and ongoing research, we are demonstrating the advantages derived from traffic filtering when placing computations other than media scheduling onto NIs, such as the downsampling of media of scientific data [42, 23], both of which use application-level header information to eliminate certain packets from a data stream. Other methods for data stream downsampling, however, like image or sensor data conversion[67], require levels of processing power not offered by NIs like the i960RD boards used in this paper. They require additional processing power, for which we are experimenting (1) with FPGA devices attached to NIs[26, 27, 30] and (2) with network processing micro-engines on IXP boards [68].

3 Insights and Experimental Evaluation

In addition to validating the feasibility of DWCS-based media scheduling on NIs, this section also demonstrates performance advantages derived from this approach.

3.1 Experimental Setup

Experiments use a typical PC-based server platform, in our case a Quad Pentium Pro server (4 X 200MHz CPUs) running Solaris 2.5.1 X86 with 128 MB of memory. Three NI cards are placed into separate PCI slots on the same bus segment. One NI hosts the scheduler and scheduler data structures. The other two cards serve as stream sources for MPEG traffic. Disks used as stream sources are directly attached to the NIs and to hosts. An MPEG segmentation program [63, 62] is used to partition an MPEG-encoded file into I, P and B frames, thereby emulating the MPEG file segmentation process in an MPEG player. The MPEG segmentation process may be run on the host CPUs or the NIs. MPEG stream producers on the host or NIs inject frames into the scheduler queues on the scheduler NI using PCI bus transfers. The scheduler picks frames based on scheduling criteria and dispatches them to the network. Client machines running MPEG players may attach to the scheduler card for MPEG stream delivery. Built-in monitoring mechanisms measure desired performance parameters at the scheduler card or at the remote client end. Remote client machines connect to the scheduler NI using a 100Mbps ethernet switched interconnect [63, 62, 21, 38, 64].

3.2 Scheduling on NIs: Basic Capabilities

Microbenchmark Definition. Microbenchmarks measure the basic performance of alternative scheduler implementations and placements. In all such benchmarks, the scheduler is started only after all frames have been written into its circular buffer, which then contains the addresses of frame descriptors. ‘Total Sched Time’ is the time to schedule all of the frames denoted by these addresses and to dispatch them to the network. ‘Avg Frame Sched Time’ is the average time to schedule a *single* frame on the network. ‘Total time w/o Scheduler’ is the cumulative time to transmit all frames on the network without the scheduler. These measurements are attained by re-routing execution in the code to a point where the address of the frame to be dispatched is readily available and does not need scheduler rules. ‘Avg Frame Time w/o Scheduler’ is the average time to transmit a single frame without the scheduler.

NI-based Scheduling is Feasible. Scheduling overhead on the i960 RD NIs is experimentally shown to be $\approx 67\mu\text{s}$. This compares favorably with the scheduling latency of the host-based scheduler reported in [63, 62] as $\approx 50\mu\text{s}$. It also corresponds to about half an ethernet frame time ($\approx 120\mu\text{s}$) on a 100Mbps link, thus indicating that it is viable to schedule multi-packet MPEG frames at network speeds. It also suggests, however, that finer-grain (e.g., per packet) scheduling would consume too much of the NI’s CPU cycles and thus, require processing resources beyond those available on this and likely, on other NIs.

NI-based Scheduling Requires Tuning. Table 1 records microbenchmarks for both a software floating point version and a fixed-point version of the DWCS scheduler. The table also shows that efficient scheduler operation on the relatively slower NIs requires some degree of tuning. For example, floating point computations are used in loss ratio computations performed in DWCS. However, like other NIs, the i960RD does not have a floating point unit. Wind River Systems (see [64]) has provided a software floating point(FP) library that may be configured into the VxWorks kernel. Measurements depicted in Table 1 show that software floating point computation results in undue scheduling overheads. We address this issue by development of a DWCS-specific fixed-point version of the library, where arguments are simply stored as fractions with numerator and denominator, with divisions implemented as shifts. This reduces computation latency by $20\mu\text{s}$, resulting in a latency of $\approx 78\mu\text{s}$ for a DWCS scheduling decision (i.e., the difference between ‘Average frame Sched time with scheduler’ and ‘Average frame Sched time without the scheduler’) for the case of fixed-point computation. The quality of scheduling, expressed in terms of parameters like delay-jitter, loss and throughput, is not affected, because the explicit representation and manipulation of numerators and denominators have the numerical accuracy required by the scheduler’s operation.

Microbenchmark	Software FP (μ Secs)		Fixed Point (μ Secs)	
	Cache Disabled	Cache Enabled	Cache Disabled	Cache Enabled
Total Sched time	19580.88	17398.56	16425.36	14295.60
Avg frame Sched time	129.67	115.20	108.48	94.60
Total time w/o Scheduler	5210.88	4776.48	4583.28	4195.68
Avg frame time w/o Scheduler	34.6	31.40	30.35	27.78

Table 1: Scheduler Microbenchmarks (Data cache effects)

Additional reductions in scheduling latency on the NI demand that the data cache on-chip be enabled¹. Table 1 shows that the presence of the data cache improves average frame scheduling time for both the software FP and the fixed point implementations of DWCS by $\approx 11\mu s$. These improvements are due to the caching of scheduler data structures, specifically, of the stream priority values and descriptor addresses, which are updated every scheduler cycle. As scheduling decisions are made on a frame-by-frame basis, data caching has the effect of reducing the ‘Total Scheduling time’ by $\approx 2182\mu s$ and $\approx 2130\mu s$, respectively, for the software FP and fixed-point implementations. Scheduler decision latency is $\approx 67\mu s$ (i.e., the difference between ‘Average frame Sched time’ and ‘Average frame time without the scheduler’) for the fixed point version.

Runtime NI Extension is Important. Data caching is one reason for using multiple NIs with each cluster node. By dedicating an NI to a specific task (i.e., scheduling vs. disk access) and thereby separating the NIs that produce media streams resident on attached disks from NIs that schedule such streams, scarce NI resources can be specialized to perform these tasks efficiently. For the i960 RD NIs, this means that the scheduler thread can benefit from data caching, without being limited by the disk driver that disables the data cache [63, 62, 21, 38, 64]. More generally, this fact indicates that the software architecture of NIs must permit the runtime *extension* of NI functionality, so that NIs can be dynamically specialized for the diverse tasks they must perform on behalf of applications running on the host nodes they are connecting[21]. A simple extension interface for an NI is described in [45]. Our ongoing work is generalizing this interface to dynamically configure the NI’s software and its attached FPGA or other specialized stream processing hardware.

Hardware Queueing is of Limited Utility. The i960 RD cards provide a number of hardware resources for device

¹Data caching has to be explicitly enabled on our i960 RD NI because the VxWorks disk driver currently supports disk accesses only with data cache disabled (the disk driver disables the data cache automatically on reboot). Therefore, for the measurements in Table 1, we first read the MPEG file from the NI-attached disk and populate the NI-resident circular buffer. After this, we enable the data cache, since further accesses to the locally attached disk are not required.

operation. These include outbound and inbound circular queues and index registers. The ‘Hardware Queues’ on the i960 RD card are a set of 1004 32 bit memory-mapped registers in local card address space. Accesses to the memory-mapped registers do not generate any external bus cycles (off-processor core). To investigate whether indexing into a circular buffer of frame descriptors may be done faster if their addresses reside in memory-mapped register space, we implemented a circular buffer where each 32 bit register holds the descriptor (with address and other attributes) of an MPEG frame. Measurements similar to Table 1 were completed. The results described in [28] indicate no additional performance advantages derived from hardware-based descriptor queues and more generally, they again demonstrate the relative paucity of NI-based hardware resources suitable for wide ranges of application-specific computations. We derive from this paucity the need for dynamic specialization and tuning of application-specific computations when they are mapped to NIs. We have had similar experiences with more modern NIs like Alteon’s gigabit ethernet boards[55] and even with ‘richer’ boards, such as Intel’s IXP 1200 router boards[52, 18, 68].

Future NIs can Schedule Media Frames for Gigabit Links. We have demonstrated that card-to-card PCI transfers may be completed without host involvement, thereby making the use of multiple NIs on a single host advantageous, and also, leaving the host CPU free to do other tasks [63, 62, 64]. We have established that it is feasible to offload scheduling functionality from the host to NIs while still meeting the frame-time requirements of MPEG frames. In general, for a scheduling discipline to schedule packets at wire-speeds and to achieve full-link utilization, decisions must be completed within ethernet frame times. Results from this section show that even on a 66 MHz i960 RD processor, the scheduler can pick winner-entities in $\approx 65\mu s$, which is within an ethernet frame time of $120\mu s$ at 100Mbps for a 1500 byte frame.

For full (max MTU) ethernet frames, the i960-based CoProcessors can handle per-frame scheduling for up to 100Mbps ethernet links. Furthermore, larger scheduling units like MPEG-I frames (each comprised of multiple ethernet frames) may be schedulable even for gigabit ethernet links.

A number of vendors are including faster and richer CoProcessors on NI boards. Alteon [55] AceNic includes two MIPS cores and the IXP1200 [52] from Intel has a StrongArm core and six hardware RISC microengines clocked at 200MHz. With substantially faster processors, it is likely that media schedulers running on these boards can meet the packet-time requirements of Ethernet frames at gigabit link speeds, as indicated by initial results described in [68] and Table 2.

Table 2 shows the results from running the packet scheduler (similar conditions as in Table 1, i.e., with data-

cache support and integer version of the packet scheduling algorithm) on a StrongArm simulator (evaluation version 1.1 from ARM) for two different configurations. Both configurations have compiled code (ARM gcc compiler flag -O1) being run directly by the simulator without any operating system (no external interrupts but with memory management). These configurations are similar to code running directly on the IXP1200 RISC microengines (without operating system support). Results from the original i960 configuration are shown in Table 2 for comparison (with VxWorks support). A StrongArm SA-110 configuration with processor and memory speeds similar to the i960 configuration yields a scheduler latency of $\approx 17.48\mu s$, which is substantially lower than the i960 scheduler latency. The third configuration is similar to the IXP1200 RISC microengines (clocked at 200MHz without any OS support) at 287 MHz and with a memory bus speed of 95.7 MHz. This shows a scheduler latency of $\approx 5.21\mu s$, thereby capable of scheduling at least 1500-byte frames at gigabit link rates. We expect to be able to lower this with careful optimization of code and compiler-assisted data placement in cache (also see [68]). We conclude from these measurements that software versions of the DWCS packet scheduler hold great promise in supporting gigabit link scheduling using modern processor chips.

Processor	Scheduler Latency(μs)
i960RD 66MHz, 4K I-cache 2K D-cache, 44MHz memory bus (with OS)	67
StrongArm SA-110 85.7MHz, 16k I-&D-cache, 28.6MHz memory bus (no OS)	17.48
StrongArm SA-110 287MHz, 16k I-&D-cache, 95.7MHz memory bus (no OS)	5.21

Table 2: Scheduler Latency for different processor configurations. cf. Ethernet Frame times for 1500 byte frames are $12\mu s(120\mu s)$ and $0.512\mu s(5.12\mu s)$ for 64 byte frames at 1Gbps (100Mbps). MPEG frames have longer allowable frame times (larger granularity).

Multiple Gigabit Links or 10Gbps Links Require Custom Packet Scheduling Hardware. Current servers already support two to three 1-Gbps interfaces, to match the backplane bandwidth of the internal PCI backplane interconnect (4.2Gbps). With the arrival of 10Gbps Infiniband and 10Gbps Ethernet hardware ([3]), the software solutions for packet scheduling explored in this paper are unlikely to scale to large numbers of packet streams. A follow-up research effort has completed construction of a custom hardware scheduler for 10Gbps links. The ShareStreams hardware prototype consists of a Xilinx Virtex I/Virtex II CoProcessor for real-time streaming under host processor systems software control[30]. While we refer the reader to [30] for details regarding the systems and hardware architecture and results attained for the Virtex II FPGA chip. Insights depicted in Table 3 demonstrate that custom hardware support is needed for scalability. With such custom hardware, even for 1024 streams, the

scheduler latency can approach 1500-byte packet-times for 10Gbps links ($1.2\mu\text{s}$):

Stream Count	Scheduler Latency (μs)	Comments
4	0.0397	meets 10Gbps packet-times
32	0.132	meets 10Gbps packet-times
1024	1.65	nearly meets 10Gbps packet-times

Table 3: Scheduler latency with Virtex II FPGA hardware

An issue for custom scheduling CoProcessors not yet addressed in this paper is that considerable input bandwidth is required to the packet scheduling hardware unit, as arrival-times for packets must be provided to the scheduler hardware unit (32-bit values) every $1.2\mu\text{s}$. PCI-66MHz (4.2Gbps peak) buses and PCIX-133MHz (8.5Gbps peak) buses do not have the backplane bandwidth to support 10Gbps links. We expect upcoming Infiniband intra-computer (system internal) interconnects to provide the internal bandwidth to support 10Gbps NIs [7]. For current Gbps NIs however, enough internal bandwidth exists with PCI and PCI-X to support multi-Gbps cross-traffic and transmission of packet-arrival times to a specialized scheduler hardware unit.

3.3 Media Scheduling and Streaming on NIs: Opportunities and Advantages

CoProcessors can affect media scheduling in ways other than filtering frames in real-time and the consequent removal of load from cluster nodes' I/O infrastructures, CPUs, and memories. We next evaluate the ability of NIs to perform more complex tasks, such as the retrieval of media content from disks directly attached to them. We also demonstrate that media scheduling on a CoProcessor is not subject to perturbation caused by system overloads. An intuitive reason for this fact is the relatively simpler nature of NIs compared to host hardware and software. This fact, coupled with the smaller, only slowly changing set of tasks a typical NI performs, makes it easier to add functionality to the NI without perturbing its other tasks or at minimum, it makes it easier to diagnose the degree of perturbation such tasks will experience. This is not true for host CPUs, even when they are multiprocessors like the quad Pentium Pro machines used in this research, when processors are dedicated to run stream schedulers, and when using all means possible to separate host-resident stream schedulers from other host tasks. Experimental results validating this claim appear later in this section.

There are multiple alternatives in media streaming from NI-attached disks. One promise of server-attached I/O CoProcessors is that simple tasks can be performed entirely independently of host CPUs, thus freeing them

for other duties and/or creating a more scalable server system comprised of the host and a number of relatively low-cost I/O CoProcessors. For high performance applications, researchers have derived benefits from the concurrency and asynchrony of execution of host vs. NI tasks, and from the low-latency access most NIs have to the actual network transceiver[45, 48, 37]. For media scheduling, we next evaluate the ability of CoProcessors to independently stream media frames to clients, once the host CPU has identified the appropriate media files and their storage or remote sources. The experiments shown below evaluate (1) whether a single NI can be both a data source and perform media scheduling, (2) how two NIs, one acting as the data source, the other performing scheduling, compare in performance to the first configuration, and finally, (3) whether the high levels of performance offered by host-resident file systems can be matched by NI-attached disks and their file system support. Option (1) is feasible for our NIs, because each i960RD card has two SCSI ports and two 100 Mbps Ethernet ports, as shown in Figure 1. Disks may be attached to the card's SCSI ports, and media may be streamed directly to the network using the card's 100 Mbps ethernet port. We also investigate Option (2), because for other NIs, we can assume the presence of peer-to-peer PCI support and an ability to change firmware to accommodate media scheduling, but disk accesses may have to be performed either via network-attached disk devices (i.e., media data is streamed to the NI via a network link) or from an intelligent disk controller [1]. Our experience is that for buses like PCI with multi-transaction timers and inter-device addressing capabilities, peer-to-peer transfers can be performed in an efficient manner. We note that both for Options (1) and (2), host CPUs like UltraSparcs and Pentium IIs are insulated from the I/O bus transfers involved through the UltraSparc Data buffer (UDB) and the PCIset, respectively.

NIs can stream data from their attached disks with performance comparable to that achieved by hosts. This statement is validated experimentally in the remainder of this section. Specifically, consider a host CPU-based scheduler that uses host filesystem buffers to store media frames, consumes I/O and system bus bandwidth, host memory and kernel/user space buffers for dispatching frames to the network. An MPEG file resident on disk must be transferred to the host's filesystem buffers by the disk controller via the I/O bus. This is shown as Path A in Figure 2. In comparison, a network interface card with an attached disk acting as a media source can use peer-to-peer PCI transfers to move data from disk, to scheduler input queues, to the network, thereby eliminating the use of host-based resources, including the host system bus. This is shown as Path B in Figure 2. Path C streams frames from a disk directly attached to the NI through to the network eliminating PCI I/O bus bandwidth consumption, host-bus bandwidth and host-memory resources in Figure 2.

We next present results from critical path benchmarks recorded for three different configurations of frame transfers. All benchmarks measure the latency of a 1000 byte frame transfer from disk to remote client (over an ethernet network) averaged over 1000 transfers. The same physical disk device is used in all experiments. The measurements in Table 4 record the latency of a single frame transfer.

NI-attached disks approximate the performance of fast host I/O systems. Consider Experiment I in Table 4, represented by Path A in Figure 2. An MPEG file on an internal system disk attached to a disk controller on the PCI bus is streamed to a remote client. The system disk is attached to a SCSI controller card on the PCI bus. The results shown in Table 4, Experiment I, show a total frame transfer time via the network of 8 ms (including disk access latency) when using the VxWorks filesystem on the Solaris host. This compares unfavorably with the results in Table 4, Experiment II, where a total of 5.4 ms is required to perform the same action for a file already resident on the NI's attached disk (Path C in Figure 2). It does indicate, however, that scalability in the number of media streams serviced by a single host can be improved by addition of low-cost I/O CoProcessors directly attached to the network. These results for Experiment I were obtained on Solaris 2.5.1 with an Intel 82557-based NI, which has the same Ethernet transceiver chipset as the i960RD CoProcessor. The system disk attached to the disk controller was used to serve frames using the Intel NI. The VxWorks filesystem is a dos-based filesystem and this was mounted on the Solaris host, in order to mitigate the effects of variations in file system performance and disk layout. Thus, the latency component common to Experiments I, II and III is the disk access time which is $\approx 4.2ms$ for a single frame (see the value '4.2disk' in Table 4).

The advantages of NI-attached disks are reduced substantially when using the faster Solaris UFS filesystem on the host. In this configuration, Experiment I experiences a disk frame latency of only 1 ms due to the larger logical block size (8K) used by UFS and its disk block caching and prefetching enabled by the host's large main memory. VxWorks does not support the UFS filesystem, so that we were unable to mount it on the NIs for Experiments II and III[13, 38], but we hypothesize that its use would improve the performance of Experiments II and III substantially if the NIs have large disk buffer caches. While this is one strong recommendation we derive from this research for NI-based media access, we also note that the host-based performance advantages gained from buffering do not extend to live media (i.e., media captured and distributed in real-time). Also the NI Coprocessor can be allowed access to the filesystem disk block buffer cache *on the host*, so that frames can be accessed directly by the Coprocessor and streamed to the network using I/O bus DMAs. The NI Coprocessor can use the host filesystem buffer cache *directly* and leverage the more efficient host filesystem for bulk data transfers.

This can help alleviate the shortcomings of the NI attached filesystem in a substantial way.

Expt	Frame Transfer Path (1000 byte frame)	Frame Transfer Time (msec)
I	Disk-Host CPU-I/O Bus-Network (no load w/ cache)	1(ufs)/8(VxWorks)
II	NI Disk-NI CPU-Network (no cache)	5.4
III	Disk-I/O Bus-NI CPU-Network (no cache)	5.415 (4.2disk+1.2net+0.015pci)

Table 4: Critical Path Benchmarks

It is advantageous to use multiple NIs, each with specialized tasks. A property of CoProcessors is that the richer their functionality, the larger the latencies of their message transfers and possibly, the smaller the total message throughput they support (see Section 3.2) [54, 37]. Our approach to this problem is to specialize NIs such that each NI only performs a limited number of tasks. In the case of media streams, this means that one NI has an attached disk, the other acts as a media scheduler and network gateway. Only small additional overheads are experienced by this NI configuration (also see Path B in Figure 2). Specifically, for a single media frame, the required peer-to-peer PCI transfers add only about $15\mu s$ to the total time of accessing a frame, scheduling it, and dispatching it to the network (we used DMA writes from card-to-card to achieve this). We achieved transfers of 66.27 MB/s along with PIO read and write latencies of $3.6\mu s$ and $3.1\mu s$ on a 32-bit 33MHz PCI bus.

The experimental results described thus far have demonstrated the viability of NI-based disk attachment and media streaming, which effectively removes the host CPU from the execution of such relatively straightforward server actions.

NI-based streaming has inherent performance advantages. To generalize our cost arguments concerning NI-based data streaming and scheduling, consider that the action of streaming frames involves (1) selecting a stream from a set of streams (sched), (2) accessing a frame from the selected stream (access), and (3) transferring the frame to the client (transfer). Stated more formally, let t_{sched} denote the time to select a stream and $t_{delivery}$ denote the time taken to deliver a single frame from the selected stream to the output link. Let t_{access} be the time to access the frame (from storage) and $t_{transfer}$ denote the transfer latency to the output link. Then:

$$t_{delivery} = t_{access} + t_{transfer}$$

and:

$$t_{stream} = t_{sched} + t_{delivery}$$

where t_{stream} is the total time for streaming a single frame.

In the case of *host-based* streaming, $t_{delivery}$ involves the relatively complex Path A shown in Figure 2, whereas for a CoProcessor, $t_{delivery}$ involves the less complex Paths B or C in Figure 2, the latter not involving any host/IO bus domain traversals. In conditions of low load and extrapolating from the measurements depicted in Table 4 (assuming the same storage source and file system capabilities for both CoProcessor and host), t_{stream} for host-based streaming is ≈ 8.050 ms (host scheduler overhead is $\approx 50\mu s$) and $\approx 5.465ms$ for CoProcessor-based streaming for a single 1000 byte frame. The advantages of CoProcessor-based streaming are exacerbated when there are host-CPU loads, as shown in Section 3.4, since t_{stream} for a host can experience significant increases even for transient loads. Further, note that t_{sched} is necessarily experienced by every frame that requires scheduling, as stream selection is done on a packet-basis and therefore, cannot be pipelined. As a result, if the term $t_{stream}^{HostLoad}$ represents t_{stream} with host load, then for a host CPU, $t_{stream}^{HostLoad} > t_{stream}$. On the other hand, for a NI CPU, host load will not affect t_{stream} , as evident in Path C in Figure 2.

3.4 Perturbation of NI- vs. Host-based Scheduling under Load

NI-based Scheduling has Reduced Delay-Jitter. The previous section has already argued that a packet scheduler running on the host (i.e., see Path A in Figure 2) uses the host’s memory, bus, and I/O bus resources. In comparison, NI-based scheduling (see Paths B or C in Figure 2) may completely avoid the consumption of I/O bus bandwidth (see Path C in Figure 2), and never uses host memory and bus bandwidth. It will therefore, be unaffected by host CPU load. This is in stark contrast to host-based scheduling, which is easily affected by the host OS’s need to run higher-level application services, where even a minimal OS installation runs system daemons and where media servers maintain meta-information in additional servers or databases. This section demonstrates experimentally the effects of CPU contention on host-based scheduling, where degradation is measured as a decrease in output bandwidth available for a stream and as an increase in its frame queueing delay. Specifically, when the DWCS scheduler receives CPU service at lower rates because of increased service load, that will lead to back-logged frames in scheduler input queues, which in turn causes missed deadlines and loss-tolerance violations. The resulting packet-dropping leads to lower scheduling quality. This is particularly important for jitter-sensitive, live-media traffic, where undue variation of the rate at which the frame/packet scheduler receives CPU services may further increase delay-jitter. Additional delay-jitter in frame output is caused by waits on congested resources, such as the multiple bus/network scheduling domains (system bus to I/O bus and then to the network) and memory

hierarchies that must be traversed by host-scheduled streams.

Experimental Demonstration of Perturbation of Host-based Scheduling in Loaded Conditions.

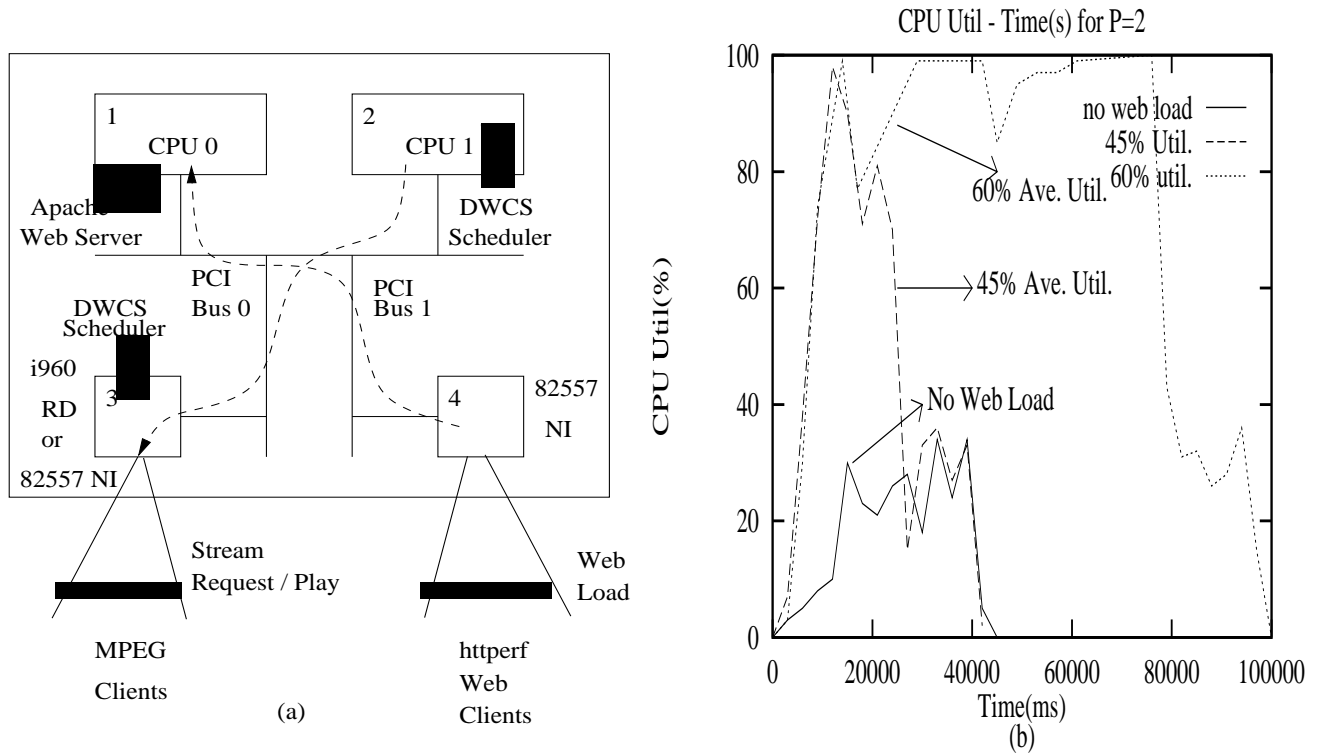


Figure 4: (a) Server Loading Architecture - Web and Media Traffic. (b) CPU Utilization Variation with Server Load.

Experimental Setup. The experimental setup consists of a Quad Pentium Pro server (4 X 200MHz CPUs) running Solaris 2.7 x86 with 128 MB of memory. This machine has two separate PCI bus segments, and we place NIs on each of the PCI bus segments, as shown in Figure 4. For host-based scheduling load experiments, Intel 82557 100Mbps transceiver-based NIs are placed in separate slots on each of the two bus segments (Components 3 and 4 in Figure 4). For NI-based scheduling load experiments, one of the Intel 82557 NIs is replaced with a i960 RD NI (component 3 in Figure 4). The machine runs the Apache web server version 1.3.12 (with a maximum of 10 server processes and starting process pool with five server processes)[6]. The web server is loaded using 'httpperf' (version 0.6)[36] run on remote Linux-based clients. Flexible specification of load from remote clients is allowed by 'httpperf', where web pages may be requested at a certain rate by a number of connections, with a user-specified ceiling on the total number of calls. The experimental infrastructure is shown in Figure 4. Placing two NIs on different PCI bus segments separates web load and stream traffic. One of the NIs (with IP address

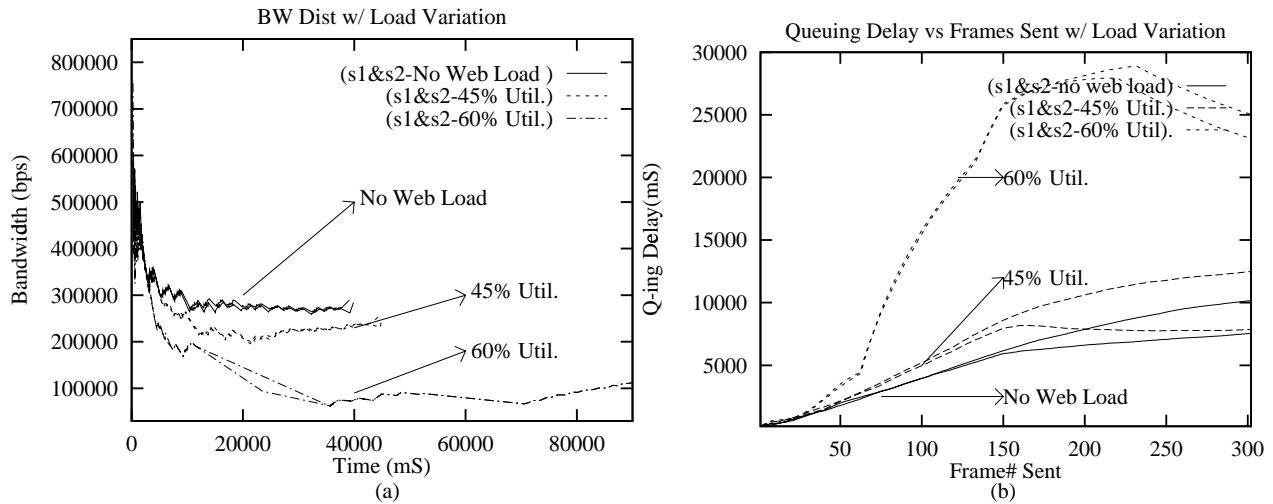


Figure 5: (a) Bandwidth Variation with Load. (b) Queuing Delay Variation with Load.

bound to the Intel 82557-based NI) is used to load the web server using ‘httperf’ client traffic while, the other NI (with a different IP address bound to the NI, 82557-based or i960 RD NI) is used to request and source stream traffic.

Moderate system loads have substantial effects on host-based scheduling. The first set of experiments involves the host-based scheduler version of the DWCS algorithm. For these experiments, two of the CPUs are brought off-line, for a total of two on-line CPUs. The Apache web server in the configuration described above is brought up and bound to the IP address of one of the NIs. The DWCS scheduler is initiated and bound to one of the host’s processors, using the ‘pbind’ Solaris facility[13]. Client requests are accepted on a separate IP address bound to a different NI. This allows ‘httperf’ web clients to connect and load the Apache Web server using a specific IP address bound to a specific NI. Similarly, MPEG clients may connect to a different IP address, again bound to a different NI, for stream delivery. This experiment involves Components 1, 2, 3 and 4, as shown in Figure 4, with Component 3 as an Intel 82557 NI. Two MPEG clients shown as Streams s1 and s2 (in Figure 5) connect to the system.

These server-load experiments demonstrate that even moderate server loads can result in substantial variations in frame bandwidth and delay. Consider, for instance, the CPU utilization depicted in Figure 4(b) (measured using Solaris’ Perfmeter facility)[13], which is the load experienced when the host-based scheduler is run with and without any load imposed by the remote web clients. With no web load, and peak utilization around 35%, with

an average utilization of 15%, the corresponding variations in bandwidth and mean queueing delay experienced for two streams s1 and s2 appear in Figures 5(a) and 5(b) (see the entries labeled ‘no web load’). In comparison, even a moderate additional host load, when applying load from web clients at the 45% utilization level, leads to noticeable differences in the observed variations in bandwidth and queueing delay. Specifically, at the 45% load level, a decrease in bandwidth to 200,000bps is seen at the 15s-20s time mark, and the bandwidth settles at only 230,000 bps (see Figure 5). The queueing delay graph in Figure 5(b) also shows the effects of loading, with frames suffering additional queueing delay of around 2s in the presence of load. A substantial web load applied at the 60% level results in what may be considered undue levels of variation. At the 60% average load level, severe degradation is seen. For instance, Figure 5 shows a decrease in bandwidth to around 100,000 bps when the CPU utilization is in the excess of 80% (see Figure 4(b)) during the period from 40s-80s. The bandwidth settles to less than 125,000 bps - half of the bandwidth seen in the absence of web server load (see Figure 5). Frames can experience excessive queueing delay in the presence of load (60% average utilization), up to three times (30,000 ms) than seen in the absence of load (10,000 ms).

NI-based scheduling is not affected by host loading. The next set of experiments involves the NI-based scheduler. For purposes of this experiment, one CPU is brought off-line (for a total of one on-line CPU), with one 82557-based Intel NI used for web server loading and a i960 RD based NI used for MPEG streaming (DWCS runs on the NI). The i960 RD NI is placed on a separate bus segment, and MPEG frames are streamed to clients by the NI-based scheduler. Loading of the web server is done using the other NI placed on a separate bus segment.

This experiment involves Components 1, 3 and 4, with Component 3 as an i960 RD NI that directly streams media content to clients. Initially, streams are played to MPEG clients in the absence of any web server load, with bandwidth variations and queueing delay measured on the NI-CPU. Next, the system is loaded using the load profile shown in Figure 4(b) (for 60% average utilization).

The measurements depicted in Figures 6(a) and 6(b) demonstrate that the NI based scheduler is immune to web server loading. In addition, no noticeable variations in bandwidth and queueing delay are experienced (see Figures 6(a) and 6(b)) for streams s1 and s2, for loaded and unloaded servers). A settling bandwidth of around 260,000 bps is observed for stream s1, which is similar to the settling bandwidth achieved by the host-CPU based scheduler in the absence of load (250,000 bps in Figure 5). These results also indicate that the i960RD NIs have sufficient ‘horse-power’ to stream scheduled media frames to clients at real-time rates.

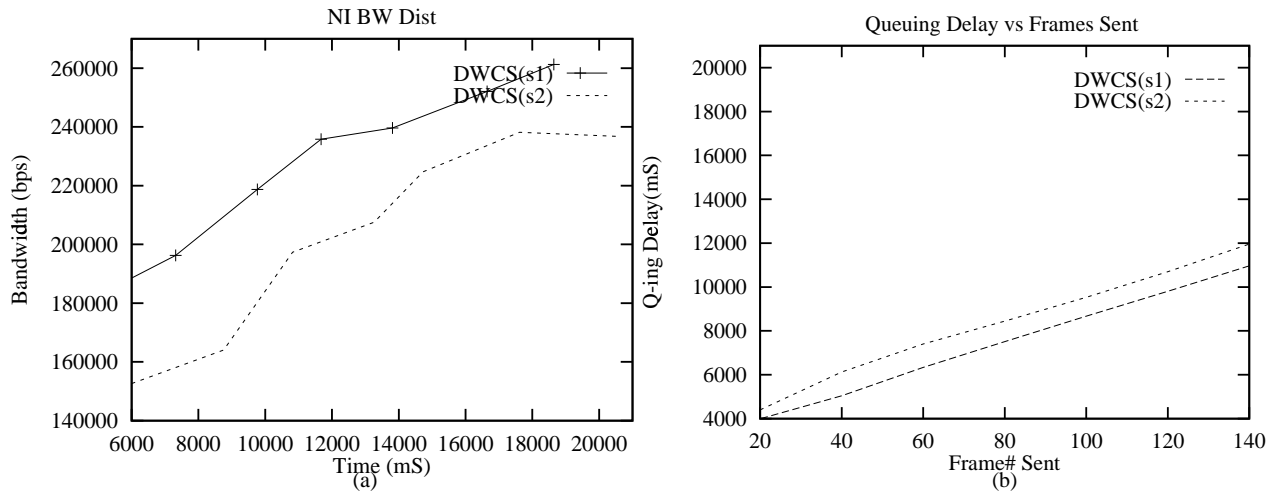


Figure 6: (a) NI Bandwidth Distribution Snapshot: Unaffected by System Load. (b) NI Queueing Delay Snapshot: Unaffected by System Load.

4 Related Work

Communication CoProcessors. A number of NI-based research projects have focused on providing low-latency message passing over cluster interconnects like ATM, Myrinet, FDDI and HIPPI[15, 40, 58, 57] by using intelligent NIs equipped with programmable CoProcessors [10, 21, 38, 45, 49]. In part, such work is motivated by the presence of programmable CoProcessors in NIs. More importantly, past research has demonstrated improved performance derived from NI reprogramming for control operations that touch multiple machines like collective communications[53], transactions[48], barriers[37, 45] and simpler cluster-wide synchronization constructs[48]. NIs have also been programmed to provide new services, like performance monitoring[34] or like the implementation of portions of protocol stacks on NIs[12].

Newer hardware developments present an interesting perspective on such research. Namely, with gigabit ethernet[55] as with the higher throughput smart port cards being developed to power the Internet's switching[32] and routing infrastructure, programmable CoProcessors (sometimes even enhanced by configurable hardware like FPGAs[30]) are becoming increasingly common. This trend coupled with the decreasing costs of processor chips suggests that future NI or I/O processor hardware will have substantial CPU cycles and memory with which additional services may be implemented. One such product is Intel's IXP 1200 router board[52], which we have begun to use as an intelligent CoProcessor for a cluster machine [68, 18].

Our work leverages prior efforts and current hardware trends by assuming that CoProcessors cannot be en-

hanced to provide all desired additional services at all times. Instead, their relative resource paucity indicates the importance of offering for CoProcessors runtime extension interfaces via which the appropriate functionality can be placed onto NIs at the right times. One such interface is described in our prior research with FORE SBA-200 (i960CA) cards [45], another in the SPINE project[17], the latter also addressing safety issues for such extensions.

The I2O industry consortium defined a specification for development of I/O hardware and software, to allow portable device driver development by defining a message-passing protocol between the host and peer I/O devices[21, 38, 64]. The focus was on relieving the host from tasks that may be offloaded to a programmable NI. Industry efforts included I2O cards for RAID storage sub-systems and off-loading TCP/IP protocol processing to the NI from the host[21, 38, 64]. While this paper uses sample I2O cards for the experimental results being presented, we do not rely on the I2O communication standard, but instead, simply use the card's CoProcessor resources. In contrast, the results presented here could benefit from the presence of improved CoProcessor-Host connectivity, as promised by the Infiniband standard [7], which offers intra-computer (within a system) and inter-computer (between systems) interconnects at 2.5Gbps, 10Gbps and 30Gbps. Improved intracomputer bandwidth using an Infiniband crossbar would help for several reasons. First, additional bandwidth from the memory subsystem to the NI would enable us to store packets in multiple places, including both the host memory subsystem or NI memory, with sufficient bandwidth available to stream packets directly from memory to output link(s). Second, storage controller channel adaptors could provide stream blocks directly to the NI using the crossbar interconnect, as the storage controller and the NI would be peers on the same Infiniband crossbar ports. The high bandwidth of an Infiniband inter-computer network would promote scalability by allowing many streams to be serviced at output link wire-speeds. In general, concerning intelligent CoProcessors in general, high performance intraprocessor interconnects like Infiniband help bridge the 'performance wall' between CPU, I/O devices, and the host processor/memory subsystem. As a result, intelligent NI-based CoProcessors would be better able to use their dedicated CPU resources to provide predictable, dedicated delivery of media streams. Similarly, if the CoProcessor is simply a scheduling engine, with media stream access and delivery performed by the host, for packet scheduling to be performed at wire speeds, the host and CoProcessor must exchange packet-arrival times of streams in every decision cycle. Even this exchange requires high aggregate bandwidth from the memory subsystem, involving 32 bits in the best case and 1023 X 32-bits (if 1023 streams drop their packets) in the worst-case for 1024 streams, every 1.2 μ s (for 10Gbps links). With RAMBUS [43] memory systems, these bandwidth requirements are only

met for the best case (4.2GBytes/sec) in current systems. Current PCI-66MHz and PCIX-133MHz cannot provide the required bandwidth, whereas Infiniband's internal crossbars should offer the capabilities suitable for scaling to future increases in network link bandwidth.

Analogous to our work with extensible NIs, there has also been research on extensible I/O CoProcessors, as with efforts like Active Disks[1] or extensible network-attached stores[31].

Media Stream Scheduling. Recent research has put substantial effort into the development of efficient scheduling algorithms for media applications. Given the presence of some underlying bandwidth reservation scheme, the DWCS algorithm has the ability to share bandwidth among competing clients in strict proportion to their deadlines and loss-tolerances. In comparison, fair share scheduling algorithms[66, 41] (some of which are now implemented in hardware[44]) attempt to allocate $1/N$ of the available bandwidth among N streams or flows. Any idle time, due to one or more flows using less than their allocated bandwidth, is divided equally among the remaining flows. This concept generalizes to weighted fairness in which bandwidth must be allocated in proportion to the *weights* associated with individual flows, but packet deadlines are not taken into account. We therefore, consider DWCS preferable for the media streams addressed by our work.

There has been significant research on the construction of scalable media servers and services, including recent work on reservation-based CPU schedulers for media applications[24]. These results demonstrate the importance of explicit scheduling to meet the demands of media applications. If DWCS performed its scheduling actions using a reservation-based CPU scheduler, it would be able to closely couple its CPU-bound packet generation and scheduling actions with the packet transmission actions required for packet streams. Similarly, DWCS could also take advantage of the stripe-based disk and machine scheduling methods advocated by some video servers [11], by using stripes as coarse-grain 'reservations' for which individual packets are scheduled to stay within the bounds defined by these reservations.

5 Conclusions and Future Work

The vision pursued by our research is one in which underlying hardware/software platforms, like the communication CoProcessors used here, are dynamically extended to better meet the needs of applications. This paper's demonstration of our vision realizes a media scheduler on an embedded NI CoProcessor using commodity hardware and software.

Insights derived from our experimental research include the following:

- *Efficient execution on standard NIs.* Even commodity NI hardware has resources sufficient for handling both the regular tasks it must perform (e.g., message receipt and sending) and certain additional tasks required by individual applications. In fact, experimental results attained with the DWCS packet scheduler for media applications demonstrate that even older, relatively slow i960-based NIs can perform a variety of such tasks, at real-time rates and at a granularity of scheduling suitable for MPEG media streams.
- *Improved scalability and predictability for NI-based streaming services.* By running media scheduling services on NIs rather than host CPUs, the host's CPU, memory, and I/O infrastructure are offloaded and in addition, substantial improvements are attained in the predictability of media streaming, measured as improvements in the delay-jitter of media streams.
- *Performance improvements due to traffic elimination.* An advantage derived from placing media scheduling onto NIs is the elimination of traffic from the host node. This fact is strengthened by the 'filtering' property of media scheduling in which losses are allowed. This property is shared by many other applications being investigated in our current research, including the selection of subsets from sensor data or from large scientific or engineering data and data downsampling or compression. The latter typically require more processing power than is currently available on commercial NIs.

To summarize, the approach to scalability for media streaming advocated by this paper has three components. The first is scalable scheduler algorithm design[63, 62] and the efficient realizations of such algorithms. The second is the appropriate distribution of media scheduling across available hardware/software resources (in this case i960 RD cards, PCI bus segments and disks), in order to separate the resources used for media scheduling from those used by general host applications. This approach also ensures that host-based programs and loading conditions do not interfere with NI-based media streaming and scheduling.

Future Work. One insight from this paper is that CoProcessors cannot be assumed capable of scheduling packets at the per-frame rates required by future communication links (e.g., 10Gbps ethernet (10GEA standard) or Infiniband). To address this issue, we have built QoS packet scheduling hardware architectures that can meet the wire-speeds of 10Gbps links. The FPGA hardware CoProcessor can meet the wire-speeds of 10Gbps links under systems software control of a peer or Host-NI CoProcessor. We are currently investigating scaling of this architecture to a large number of streams by allowing stream state to be multiplexed on the same FPGA hardware substrate [60, 26, 27, 30]. In addition, we are studying how to 'divide' scheduling functionality or how to simplify it so

that per frame scheduling may be performed with the RISC-based microengines resident on typical next generation CoProcessors, such as Intel's IXP1200[68]. Finally, beyond packet scheduling, we are investigating suitable extension architectures for NIs, so that application-specific functionality is easily mapped to them, whenever applications can benefit from such NI-based support. One set of extensions already under development by our group concerns the reliability and scalability of transaction services implemented on cluster machines[19, 20].

Acknowledgements. This work was supported in part by the Department of Energy under its NGI program, by the NSF, and by hardware/software donations from Intel Corporation and WindRiver Systems. We would like to thank the reviewers for their many invaluable suggestions. Our thanks also go to other contributors to this research, including Prof. Sudhakar Yalamanchili, Prof. Ken Mackenzie and MS students S. Manni and S. Roy.

References

- [1] A. Acharya, M. Uysal, and J. Saltz. Active disks. In *International Conference on the Architectural Support for Programming Languages and Operating Systems (ASPLOS98)*, 1998.
- [2] Akamai and FreeFlow content management. <http://www.akamai.com>.
- [3] 10 Gigabit Ethernet Alliance. <http://www.10gea.org>.
- [4] Kevin Almeroth and Mostafa Ammar. A scalable, interactive video-on-demand service using multicast communication. In *Proceedings of the International Conference on Computer Communication Networks*, San Francisco, California, September 1994.
- [5] Thomas E. Anderson, David E. Culler, David A. Patterson, and the NOW Team. A Case for Networks of Workstations: NOW. *IEEE Micro*, Feb. 1995.
- [6] Apache http Server Project Apache Software Foundation. <http://www.apache.org/httpd.html>.
- [7] Infiniband Trade Association. <http://www.infinibandta.org>.
- [8] Brian N. Bershad, Stefan Savage, Przemyslaw Paradyk, Emin Gum Sirer, Marc Fiuczynski, and Becker Eggers Chambers. Extensibility, safety, and performance in the SPIN operating system. In *Proceedings of the 15th ACM Symposium on Operating Systems Principles*, pages 267–284, Copper Mountain, Colorado, December 1995.
- [9] Ben Blake. *A Fast, Efficient Scheduling Framework for Parallel Computing Systems*. PhD thesis, Department of Computer and Information Science, The Ohio State University, Dec. 1989.
- [10] Nanette J. Boden, Danny Cohen, Robert E. Felderman, Alan E. Kulawik, Charles L. Seitz, Jakov N. Seizovic, and Wen-King Su. Myrinet – A Gigabit-per-Second Local-Area Network. *IEEE MICRO*, Feb. 1995.

- [11] William J. Bolosky, Robert P. Fitzgerald, and John R. Douceur. Distributed schedule management in the tiger video fileserver. In *Sixteenth ACM Symposium on Operating System Principles*, volume 31, pages 212–223. ACM, December 1997.
- [12] C.Keppitiyagama and A. Wagner et al. Asynchronous mpi messaging on myrinet. In *Int'l Parallel and Distributed Processing Symposium (IPDPS '01) San Francisco*, April 2001.
- [13] Solaris On-Line Documentation. <http://www.docs.sun.com>.
- [14] Dawson R. Engler, M. Frans Kaashoek, and James O'Toole Jr. Exokernel: An operating system architecture for application-level resource management. In *Proceedings of the 15th Symposium on Operating System Principles*, December 1995.
- [15] Edward W. Felten, Richard D. Alpert, Angelos Bilas, Matthias A. Blumrich, Douglas W. Clark, Stefanos Damianakis, Cezary Dubnicki, Liviu Iftode, and Kai Li. Early Experience with Message-Passing on the SHRIMP Multicomputer. *Proceedings of the 23rd International Symposium on Computer Architecture*, May 1996.
- [16] D. Ferrari, A. Banerjea, and H. Zhang. Network support for multimedia. a discussion of the tenet approach. *TR-92-072 UC Berkeley, Department of Computer Science*, 1992.
- [17] Marc E. Fiuczynski, Brian N. Bershad, R.P. Martin, and D.E. Culler. SPINE - An Operating System for Intelligent Network Adapters. (TR-98-08-01), Aug. 1998.
- [18] Ada Gavrilovska, Karsten Schwan, Austen McDonald, and Ken Mackenzie. Stream handlers: Application-specific message services on attached network processors. In *Proceedings of the 10th IEEE Conference on High-Performance Interconnects (HotI-02), Stanford University, California*, Aug 21-23, 2002.
- [19] Ada Gavrilovska, Karsten Schwan, and Van Oleson. Adaptable mirroring on clusters. In *Tenth International Conference on High Performance Distributed Computing(HPDC-10), San Francisco, California.*, August 2001.
- [20] Ada Gavrilovska, Karsten Schwan, and Van Oleson. Practical approach to zero downtime in an operational information system. In *Proceedings of the 22nd IEEE International Symposium on Distributed Computing Systems*. IEEE, 2002.
- [21] I₂O Special Interest Group. www.i2osig.org/architecture/techback98.html.
- [22] Intel. *IQ80960Rx Evaluation Platform Board Manual*, March 1997.
- [23] C. Isert and K. Schwan. Acds: Adapting computational data streams for high performance. In *Proceedings of International Parallel and Distributed Processing Symposium (IPDPS 2000)*, 2000.
- [24] Michael B. Jones, Daniela Rosu, and Marcel-Catalan Rosu. Cpu reservations and time constraints: Efficient, predictable scheduling of independent activities. In *Sixteenth ACM Symposium on Operating System Principles*, volume 31, pages 198–211. ACM, December 1997.
- [25] Charles Krasic and Jonathan Walpole. Qos scalability for streamed media delivery. Technical Report CSE-99-011 Oregon Graduate Institute, Dept. of Computer Science, 17, 1999.
- [26] Raj Krishnamurthy and K. Azad et al. The georgia tech asan approach. In *IEEE International Symposium on High Performance Computer Architecture (HPCA-7). Work-in-Progress Session. Monterrey, Mexico*, January 2001.

- [27] Raj Krishnamurthy and Sudhakar Yalamanchili et al. Architecture and hardware for scheduling of gigabit packet streams. In *IEEE International Symposium on High Performance Computer Architecture (HPCA-7). Work-in-Progress Session. Monterrey, Mexico, January 2001.*
- [28] Raj Krishnamurthy, Karsten Schwan, Richard West, and Marcel Rosu. A network coprocessor-based approach to scalable media streaming in servers. Technical Report GIT-CC-00-03, Georgia Institute of Technology, 2000.
- [29] Raj Krishnamurthy, Karsten Schwan, Richard West, and Marcel Rosu. A network co-processor-based approach to scalable media streaming in servers. In *Proceedings of the International Conference on Parallel Processing (ICPP-2000)*. International Association for Computers and Communication(IACC), Aug 21-24, 2000.
- [30] Raj Krishnamurthy, Sudhakar Yalamanchili, Karsten Schwan, and Richard West. Architecture and hardware for scheduling gigabit packet streams. In *in the Proceedings of the 10th IEEE Conference on High-Performance Interconnects (HotI-02), Stanford University, California, Aug 21-23, 2002.*
- [31] E. Lee and C. Thekkath. Petal: Distributed virtual disks. In *International Conference on the Architectural Support for Programming Languages and Operating(ASPLOS96)*, 1996.
- [32] John Lockwood, Jon Turner, and David Taylor. Field programmable port extender(fpx) for distributed routing and queuing. In *ACM International Symposium on Field Programmable Gate Arrays (FPGA'2000), Monterey, CA, pp. 137-144, February 2000.*
- [33] Gordon Mair. Telepresence - the technology and its economic and social implications. In *Proceedings of the IEEE International Symposium on Technology and Society*. IEEE, 1997.
- [34] M. Martonosi, D. Clark, and M. Mesarina. The shrimp hardware performance monitor: Design and applications. In *Proc. 1996 SIGMETRICS Symposium on Parallel and Distributed Tools*, 1996.
- [35] C.W. Mercer, S. Savage, and H. Tokuda. Processor capacity reservation for multimedia operating systems. In *IEEE International Conference on Multimedia Computing and Systems*. IEEE, May 1994.
- [36] David Mosberger and Tai Jin. httpperf – a tool for measuring web server performance. In *Proceedings of the 1998 Workshop on Internet Server Performance , held in conjunction with Sigmetrics 1998*, June 1998.
- [37] J. Nieplocha and J. Ju et al. One-sided communication on myrinet-based smp clusters using the gm message-passing library. In *Workshop on Communication Architectures in Clusters held in Conjunction with Int'l Parallel and Distributed Processing Symposium (IPDPS '01) San Francisco, April 23-27, 2001*, April 2001.
- [38] *I₂O Intel Page*. <http://www.developer.intel.com/iio>.
- [39] V. Pai, P. Druschel, and W. Zwaenepoel. Io-lite: A unified buffering and caching system. In *3rd Symposium on Operating Systems Design and Implementation (OSDI '99) Proceedings.*, 1999.
- [40] Scott Pakin, Mario Laura, and Andrew Chien. High Performance Messaging on Workstations: Illinois Fast Messages (FM) for Myrinet. *Supercomputing*, Dec. 1995.
- [41] Xingang Guo Pawan Goyal and Harrick M. Vin. A hierarchical cpu scheduler for multimedia operating systems. In *2nd Symposium on Operating Systems Design and Implementation*, pages 107–121. USENIX, 1996.

- [42] Beth Plale and Karsten Schwan. dquob: Managing large data flows by dynamic embedded queries. In *Proceedings of the IEEE International Symposium on High Performance Distributed Computing (HPDC 2000)*. IEEE, 2000.
- [43] RAMBUS. <http://www.rambus.com>.
- [44] Jennifer L. Rexford, Albert G. Greenberg, and Flavio G. Bonomi. Hardware-efficient fair queueing architectures for high-speed networks. In *INFOCOMM'96*, pages 638–646. IEEE, March 1996.
- [45] Marcel-Cătălin Roșu, Karsten Schwan, and Richard Fujimoto. Supporting Parallel Applications on Clusters of Workstations: The Intelligent Network Interface Approach. *Proceedings of the 6th IEEE International Symposium on High Performance Distributed Computing*, Aug. 1997.
- [46] Daniela Rosu, Karsten Schwan, and Sudhakar Yalamanchili. FARA - a framework for adaptive resource allocation in complex real-time systems. In *Proceedings of the 4th IEEE Real-Time Technology and Applications Symposium (RTAS)*, Denver, USA, June 1998.
- [47] Daniela Rosu, Karsten Schwan, Sudhakar Yalamanchili, and Rakesh Jha. On Adaptive Resource Allocation for Complex Real-Time Applications. *18th IEEE Real-Time Systems Symposium*, Dec., 1997.
- [48] Marcel-Catalin Rosu and Karsten Schwan. Sender coordination in the distributed virtual communication machine. In *Proceedings of the 7th IEEE International Symposium on High Performance Distributed Computing (HPDC '98)*. IEEE, 1998.
- [49] Marcel-Catalin Rosu, Karsten Schwan, and Richard Fujimoto. Supporting Parallel Applications on Clusters of Workstations: The *Virtual Communication Machine*-based Architecture. *Cluster Computing*, 1:1–17, Jan. 1998.
- [50] Y. Saito, B. Bershad, and H. Levy. Availability and performance in porcupine: A highly scalable internet mail service. In *Proceedings of the 17th ACM Symposium on Operating Systems Principles*, Charleston, SC, December 1999.
- [51] Karsten Schwan and Hongyi Zhou. Dynamic scheduling of hard real-time tasks and real-time threads. *IEEE Transactions on Software Engineering*, 18(8):736–748, August 1992.
- [52] Intel IXP 1200 Web Site. <http://www.intel.com/design/network/products/npfamily/index.htm>.
- [53] C. Stunkel, D. Shea, B. Abali, M. Atkins, C. Bender, D. Grice, P. Hochschild, D. Joseph, B. Nathanson, R. Swetz, R. Stucke, M. Tsao, and P. Varker. The sp2 communication subsystem. In *Tech. report, IBM Thomas J. Watson Research Center, Yorktown Heights, NY, August 1994*. WWW URL is <http://ibm.tc.cornell.edu>, 1994.
- [54] R. Swan, S. Fuller, D. Siewiorek, and C. modular. multi-microprocessor. In *1977 National Computer Conference, volume 46, pages 637–644, 1977, 1977*.
- [55] Alteon Web Systems. <http://www.alteonwebsystems.com>.
- [56] Trivedi, Hall, Kogut, and Roche. Web-based teleautonomy and telepresence. In *SPIE Optical Science and Technology Conference, 2000*. SPIE, 2000.
- [57] Thorsten von Eicken, Anindya Basu, Vineet Buch, and Werner Vogels. U-Net: A User-Level Network Interface for Parallel and Distributed Computing. *Proceedings of the 15th ACM Symposium on Operating Systems Principles*, Dec. 1995.

- [58] Thorsten von Eicken, David E. Culler, Seth Copen Goldstein, and Klaus Erik Schauser. Active messages: A mechanism for integrated communication and computation. *Proceedings of the 19th International Symposium on Computer Architecture*, May 1992.
- [59] Jonathan Walpole, Rainer Koster, Shanwei Chen, Crispin Cowan, David Maier, Dylan McNamee, Calton Pu, David Steere, and Liujin Yu. "a player for adaptive mpeg video streaming over the internet". In *Proceedings 26th Applied Imagery Pattern Recognition Workshop AIPR-9*, Washington DC, October 1997.
- [60] Richard West, Raj Krishnamurthy, William Norton, Karsten Schwan, Sudhakar Yalamanchili, Marcel Rosu, and Sarat Chandra. Quic: A quality of service network interface layer for communication in nows. In *Proceedings of the Heterogeneous Computing Workshop, in conjunction with IPPS/SPDP*, San Juan, Puerto Rico, April 1999.
- [61] Richard West and Christian Poellabauer. Analysis of a window-constrained scheduler for real-time and best-effort packet streams. In *Proceedings of the 21st IEEE International Symposium on Real-time Systems (RTSS '00)*. IEEE, 2000.
- [62] Richard West and Karsten Schwan. Dynamic window-constrained scheduling for multimedia applications. In *6th International Conference on Multimedia Computing and Systems, ICMCS'99*. IEEE, June 1999. Also available as a Technical Report: GIT-CC-98-18, Georgia Institute of Technology.
- [63] Richard West, Karsten Schwan, and Christian Poellabauer. Scalable scheduling support for loss and delay constrained media streams. Technical Report GIT-CC-98-29, Georgia Institute of Technology, 1998.
- [64] WindRiver Systems. *VxWorks Reference Manual*, 1 edition, February 1997.
- [65] Xilinx. <http://www.xilinx.com>.
- [66] Hui Zhang and Srinivasav Keshav. Comparison of rate-based service disciplines. In *Proceedings of ACM SIGCOMM*, pages 113–121. ACM, August 1991.
- [67] D. Zhou and K. Schwan. Adaptation and specialization for high performance mobile agents. In *Usenix COOTS99*, 1999.
- [68] Xiatong Zhuang, Weidong Shi, Indrani Paul, and Karsten Schwan. On the efficient implementation of the dwcs packet scheduling algorithm on ixp1200 network processors. In *Proceedings of the IEEE International Symposium on Multimedia Networks and Systems*. IEEE, 2002.