

Migrating Sockets for Networking with Quality of Service Guarantees*

David K.Y. Yau and Simon S. Lam
Department of Computer Sciences
The University of Texas at Austin
Austin, Texas 78712-1188
{yau,lam}@cs.utexas.edu

Abstract

Migrating Sockets is the protocol processing component of an end system architecture designed for networking with QoS guarantees. The architecture provides (1) adaptive rate-controlled scheduling of protocol threads in Migrating Sockets, (2) rate-based flow control for reserved rate connections in future integrated services networks, and (3) a *constant overhead* active demultiplexing mechanism. Migrating Sockets achieves its efficiency by allowing user applications to manage a network endpoint with minimal system intervention, providing user level protocols read-only access to routing information in a “well-known” shared memory region, and integrating efficient kernel level support we previously built. It is backward compatible with Unix semantics and Berkeley sockets, and has been used to implement Internet protocols such as TCP, UDP and IP (including IP multicast). We also show that active demultiplexing supported by Migrating Sockets can be *transparently* enabled in wide-area TCP/IP internetworking (although it is not restricted to TCP/IP). We have an implementation of Migrating Sockets in Solaris 2.5. We discuss our implementation experience, and present performance results of our system running on the Ultra-1, SPARC 10 and SPARC 20 architectures.

1 Introduction

It is increasingly important for end systems to provide support for networking with quality of service (QoS) guarantees. This trend is in part due to the emergence of continuous media (such as video and audio) applications having real-time constraints. Such work on end system support complements recent research on integrated services networks. QoS guarantees provided by network level packet scheduling and admission control through Internet resource reservation protocols being developed [1, 5, 16] can thus be extended to the ultimate endpoints of an end-to-end communication, namely applications running in user space of general purpose operating systems.

End system support for networking with QoS guarantees is a challenging problem. To meet various timing constraints, user processes must be given guaranteed access to diverse system resources, including time-shared resources such as CPU and network interface, and space-shared resources such as memory buffers. Moreover, the run time environment for protocol processing, which provides such services as timer management, buffer management and demultiplexing table lookup, should be designed to support predictable performance.

Besides QoS guarantees, recent proliferation of heterogeneous net-

working technologies and user application requirements will make customized development and flexible deployment of network protocols highly desirable. Protocol implementation at user level can help achieve these goals. With fault containment in user processes and the availability of sophisticated tools for developing user level code, the cost of protocol development and experimentation will go down, and the lead time to deployment of protocols in a production environment will be reduced [12]. Moreover, without the need to configure and load protocols into kernel space, user applications can be given access to a wider choice of protocol stacks and choose the one that is most suitable for their particular needs.

1.1 Our contributions

We discuss our experience implementing Migrating Sockets as a framework for user level protocols that can run with guaranteed progress. In addition, we highlight several novel ideas implemented in Migrating Sockets. First, we discuss how Migrating Sockets can minimize *hidden scheduling* in protocol processing. This allows protocol threads to run with progress guarantees within the context of adaptive rate-controlled (ARC) scheduling. Second, we present the use of *delayed caching* and *recall-on-access* for socket migration that improves caching performance for the widely used concurrent server model in client/server programming. Third, we show how the use of a “well-known” shared memory region can enable extremely efficient sharing of routing information between network and higher layer protocols in a user level implementation model. Fourth, Migrating Sockets has been integrated with operating system techniques that minimize data copying and system call overhead for network communication [14]. Lastly, “exclusive packet receiver” information exported by Migrating Sockets enables a constant overhead packet demultiplexing mechanism called *active demultiplexing*. By eliminating table search, active demultiplexing is highly efficient and is suitable for networking with QoS guarantees.

1.2 Organization of this paper

The balance of this paper is as follows. In the following subsection, we discuss related work. In section 2, we give an overview of our end system support architecture for networking with QoS guarantees. We show how Migrating Sockets works together with other system components to provide QoS guarantees in an end-to-end path of network communication. Section 3 introduces Migrating Sockets as a user level protocol implementation framework. The framework is flexible, efficient and backward compatible with Unix semantics and Berkeley sockets. Section 4 describes several novel aspects of Migrating Sockets. ARC scheduling of protocol threads is presented in section 5. Active demultiplexing is described in section 6. We present experimental results on the performance of our current system in section 7, and conclude in section 8.

Research supported in part by National Science Foundation under grant no. NCR-9506048, an equipment grant from AT&T Foundation, and an IBM graduate fellowship. David K.Y. Yau is now with the Department of Computer Sciences, Purdue University, West Lafayette, IN 47907 (email: yau@cs.purdue.edu).

1.3 Related work

Real-time upcalls were proposed in [6] to achieve QoS guarantees in protocol processing. While the approach is an interesting alternative to real-time threads, it is specifically designed for periodic protocol processing and appears to be less general than our approach. Protocol processing with predictable performance has also been investigated in the context of Real-Time Mach [10] based upon design principles for CPU scheduling different from ours.

There has been growing interest in user level protocol implementation in recent years [8, 12]. A user level TCP implementation on top of the Jetstream high-speed network interface is described in [3]. U-Net [13] integrates interface firmware with host software in a design that provides user level access to a network without kernel intervention. These works target high performance on the send/receive path, without paying much attention to the issues of connection management, routing management, and the semantics of sharing network endpoints. Performance benefits of Integrated Layer Processing in user level protocols are evaluated in [2]. Migrating Sockets similarly enables the integration of protocol functions at the presentation, session, transport and network layers.

2 Architectural Overview

Our end system architecture for networking with QoS guarantees has the following major components: (1) ARC scheduling for time-shared resources in an end system, such as CPU and network interface, (2) a migrating sockets framework for user level protocols that minimizes hidden scheduling in protocol processing, and (3) a constant overhead packet demultiplexing mechanism suitable for wide area internetworking. Figure 1 illustrates the architecture. A packet path from the sender system on the left to the receiver system on the right is shown.

Migrating Sockets (see section 3) takes its name because the state and management right of a network endpoint can move between a network server and client processes. With Migrating Sockets, performance critical protocol services such as send and receive are accessed as a user level library linked with applications. Send side protocol code is accessed in usual application threads of control. In addition, user processes have protocol threads for network receive and timer processing. From a QoS perspective, Migrating Sockets has the advantage of minimizing “hidden” scheduling.¹ For example, the role of network interrupt handlers in our system is only to deliver packets to a set of destination processes. The bulk of receive side processing is done in the context of a user thread of control.

With hidden scheduling minimized, user applications can more easily determine and negotiate an appropriate *rate of progress* with an end system, such that their real-time constraints can be met. In our architecture, progress requirements are specified with two parameters: a reserved *rate* (between 0 and 1) and a time interval known as *period* (in μs). Based on the progress requirements of all threads in the system, an ARC CPU scheduler can perform admission control and provide *conditional* progress guarantees to threads. The ARC scheduler in [15] provides the following progress guarantee: a “punctual” thread with rate r and period p is guaranteed at least $k r p$ CPU time over time interval $k p$, for $k = 1, 2, \dots$

On the send side, ARC scheduling enables applications to respond to media events and generate network packets “in time.” These packets then enter a network connection, which may have a reserved

¹Hidden scheduling occurs when protocol processing is done in the context of interrupt handling or “background” threads of control that do not belong to a user process.

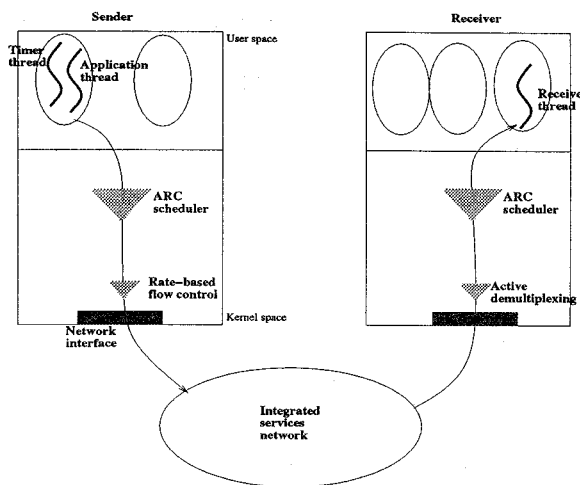


Figure 1: End system architecture for networking with QoS guarantees.

rate negotiated with a future integrated services network. If the network connection is shared by multiple processes, it is possible for a sudden burst of packets by one process to block out access to the network connection for an extended period of time [14], thereby jeopardizing the bandwidth requirements of other processes. This problem is especially pronounced if the connection has a moderate or low reserved rate. To solve the problem, an end system should provide rate-based flow control to reserved-rate network connections. In our proposal [14], flow control is enforced by a *lightweight kernel thread*. The approach is quite flexible in that different flow control policies can be provided by different loadable kernel modules.

On the receive side, packet arrivals to a network interface are processed by the interrupt handler of the interface. Kernel level code must then demultiplex the packets to their destination processes. Traditionally, such demultiplexing is performed by packet filters [4, 9]. Our system makes use of packet filters, but, in addition, can exploit “exclusive packet receiver” information exported by Migrating Sockets to perform *active* packet demultiplexing. A network endpoint that is an exclusive packet receiver has the property that packets destined for it should not be delivered to any other endpoint in the system.

In active demultiplexing, an exclusive packet receiver *advertises* to a peer sender an OS handle for packet delivery. On learning the advertisement, the sender *encloses* this OS handle in packets it sends to the receiver. The kernel demultiplexing code in the receive system can then make use of the handle to deliver packets directly to the receiver, without table searching. For safety reasons, the receive kernel checks that a handle is indeed associated with an exclusive packet receiver and ensures the “freshness” of the handle by using a *nonce* included with the handle. From a QoS perspective, active demultiplexing is desirable since it is a constant overhead mechanism, contributing to predictable performance.

In the receive end system, demultiplexed network packets may cause their receiver processes to be scheduled. An ARC CPU scheduler in the receive system enables such processes to respond to the packet arrivals “in time.”

3 Migrating Sockets

In choosing an application programming interface (API) for our system, one of our goals is that the API should allow us to run existing and future Unix multimedia applications (such as the mbone suite of teleconferencing tools) with minimal modifications. We therefore decided to maintain backward compatibility with the widely used Berkeley socket interface. In concept, our migrating sockets framework draws upon previous experience in user level protocol implementation [8, 12]. In what follows, we give an overview of Migrating Sockets. Several novel ideas implemented in the framework are presented in section 4.

3.1 Endpoint management

Figure 2 gives an overview of the protocol implementation model of Migrating Sockets. The model uses a network server process for “boundary” protocol operations like opening and closing network connections. User applications are client processes in the model. They register themselves with the network server. When a client process opens a socket, it creates a local data structure for the socket and contacts the network server with an RPC call. The network server likewise creates a data structure for the socket, and performs most of the “real” work for socket creation. Apart from using the `socket()` call, client processes can *implicitly* gain access to socket descriptors, such as after a fork. In the case of a fork, the network server copies socket descriptors from the parent process to the child process and uses callback RPC to ask the child process to create data structures for the sockets inherited. In summary, a client process knows about all the sockets for which it holds a socket descriptor, whereas the network server knows about all the sockets in an end system.

Since the network server has global knowledge of all the network endpoints in a system, it is suitable for the tasks of connection management. For example, it will be able to enforce uniqueness of network connections requested by different client processes. However, it is expensive to go through the network server for every socket operation. This is especially true of operations on the “performance critical” path, namely sending and receiving network data. In Migrating Sockets, after a network connection has been established, the state and management right of a socket can be *cached* to the client process holding an exclusive socket descriptor for the socket. Caching involves transferring the current state of the socket (including any data buffered for send and receive, and any protocol state associated with the socket) from the network server to the client process. A client process can then send to and receive from the cached socket without going through the network server.

For connectionless protocols, such as UDP, a socket can be cached after the `bind()` system call, which fills in a local address and port number for the socket. For connection oriented protocols, such as TCP, a socket can be cached after the `accept()` system call, which fills in both local and remote addresses and port numbers for the socket.

A socket continues to be cached until a condition unfavorable to caching occurs, at which time the socket is *flushed* to the network server. This involves transferring the state and management right of the socket from the client process back to the network server. As an example, flushing is required before a fork, since cached access is incompatible with the semantics of sharing network endpoints. It is also required before a close. After flushing, the network endpoint may continue to exist within the network server. Therefore, flushing before a close takes care of the requirement of certain protocols (such as TCP) that the lifetime of a network endpoint may exceed

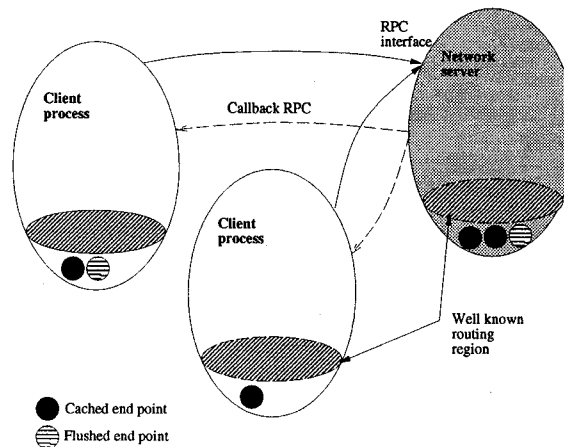


Figure 2: The protocol implementation model of Migrating Sockets.

the lifetime of the process having access to the endpoint. All operations on a flushed socket go through the network server, using an RPC interface.

Whereas shared access to a socket can prevent the socket from being cached, caching may be re-enabled by the `close()` operation. The situation occurs if, after a close by some process, another process becomes the only one holding a socket descriptor to the socket in question. When that happens, our system caches the socket to the latter process.

Besides the cache and flush operations, our system supports a third operation for socket migration known as *recall*. The operation is a callback RPC for the network server to ask a client process to transfer the state and management right of a cached socket back to the network server. It is needed, for example, in the optimization described in section 4.2.

3.2 Implementation considerations

Our implementation of Migrating Sockets leverages protocol code from 4.4 BSD. However, parts of the runtime support system have been rewritten. First, we replaced BSD mbuf buffer management by *message blocks* similar to those used in SVR4 streams. This is because mbuf has been found to treat small and large messages non-uniformly and hence exhibit undesirable performance idiosyncrasies [7]. Moreover, message blocks can very naturally handle both normal data buffers and *network buffers* (see section 4.4) supported in our system (using the `esballoc()` library call).

Second, we implemented a timer management interface for timer activities. Unlike 4.4 BSD, timer processing is driven by a timer thread of control.

Third, it is not feasible to protect critical code sections by raising interrupt level in a user level implementation. Instead, we make use of mutex locks and condition variables for mutual exclusion and condition synchronization. The current locking granularity is quite coarse. To illustrate, it is often convenient to think of protocol processing as consisting of an “upper” and a “bottom” half. The upper half is driven by protocol send activities, while the bottom half is driven by packets received from the network. Figure 3 shows the multi-threaded structure of a typical client process accessing cached sockets. A timer thread and application threads with send side protocol code run in the upper half. Protocol receive threads run in the bottom half. Notice that we use a thread of control for receiving

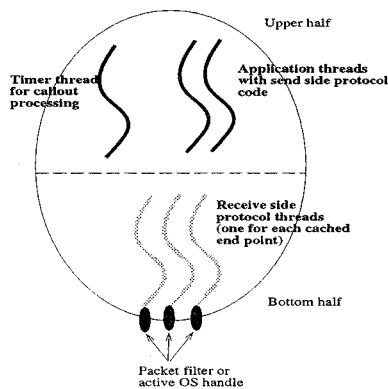


Figure 3: Multi-threaded process structure for applications accessing cached sockets.

from each cached endpoint.

In our locking model, a mutex lock `syslock` protects system data structures that are not modified by bottom half threads. The purpose of the lock is then to synchronize access by multiple upper half threads. For example, an upper half thread opening a socket and another one doing a close may both be trying to modify a set of file descriptors associated with a process. Locking is required to serialize the two operations. Another mutex lock `intrlock` protects system data structures that may be modified by bottom half threads. Its purpose is to synchronize access by bottom half threads, as well as between the upper and bottom halves. For example, a bottom half thread may be trying to append data to a receive socket buffer, while an upper half thread may be trying to remove data from it.

Upper half protocol threads normally acquire `syslock` before `intrlock`. Consistent use of this locking order avoids deadlocks between threads in acquiring the two locks. There are situations, however, in which the normal locking order is not followed. For these situations, the code fragment shown in Figure 4 is used to avoid deadlocks. Specifically, an upper half thread is trying to read from a socket. The thread first acquires `syslock` and then `intrlock` before it checks the receive socket buffer. If it finds no data available for reading, the thread blocks on a condition `readable` protected by `intrlock` (`intrlock` will be automatically released when blocking occurs inside `cond.wait()`). Before the thread calls `cond.wait()`, however, it explicitly releases `syslock` to allow other upper half threads access to data structures protected by `syslock`.

When the thread wakes up because data have arrived, it will have acquired `intrlock` while also needing to re-acquire `syslock`. In the figure, the call `mutex_trylock(syslock)` tries to acquire `syslock` but returns a failure condition of 0 if the lock cannot be acquired. Notice that if `mutex_trylock()` succeeded, the thread will have acquired both `syslock` and `intrlock`. If `mutex_trylock()` failed, however, the thread first releases `intrlock` and then tries again to acquire both locks in the normal order of `syslock` before `intrlock`, thus avoiding the possibility of deadlock with another thread.

4 Novel Aspects of Migrating Sockets

Our migrating sockets framework exports information on *exclusive packet receivers*, which has the following meaning: If a network endpoint is an exclusive packet receiver, then packets destined for it should not be delivered to any other endpoint in the system. Notice that in Migrating Sockets, network connections are established

Algorithm LOCKING

```
mutex_lock(syslock);
...
mutex_lock(intrlock);
...
while (no data available for reading) {
    mutex_unlock(syslock);
    cond_wait(readable, intrlock);
    if (mutex_trylock(syslock) == 0) {
        mutex_unlock(intrlock);
        mutex_lock(syslock);
        mutex_lock(intrlock);
    }
}
...
}
```

Figure 4: Locking algorithm for deadlock avoidance.

through the network server, which knows about all the existing network endpoints in the system. Hence, when the network server allows a socket to be cached, it knows whether the socket being cached is an exclusive packet receiver or not.

Information on exclusive packet receivers can be used to reduce the search time for matching packets with packet filters, since a match with the filter of such a receiver means that further matching would be unnecessary. Moreover, as we will show in section 6, the information enables a constant overhead packet demultiplexing mechanism known as *active demultiplexing*.

We now elaborate on several other aspects of our system that are novel or pertaining to the provision of QoS guarantees. They are (1) minimizing hidden scheduling in protocol processing, (2) caching optimization for the concurrent server programming model, (3) sharing of routing information between network and higher level protocols using a “well-known” shared memory region, and (4) a kernel/user interface that provides user level protocol code with access to efficient kernel level support [14] through Unix file descriptors.

4.1 Minimizing hidden scheduling

Our experience [15] has been that it is difficult to provide QoS guarantees in certain protocol implementation frameworks. In SVR4 streams, for example, network send and receive can take place in *service routines* run by “background” system threads of control. In BSD Unix, a single system timeout invocation has to handle outstanding timer activities of all the network endpoints in the system. The main problem of these background system services is that there is no easy way to determine suitable reserved rates of progress for the system services, such that the real time constraints of user applications can be met.

Aside from the use of background system services, traditional kernel level protocols perform entire receive side protocol processing in the context of interrupt handling. From a QoS perspective, it is similarly difficult to control the rate of progress of interrupt handling code.

Migrating Sockets reduces the use of such hidden scheduling for cached sockets. First, each user process has a dedicated timer thread that handles timer events only for network endpoints local to the

Algorithm CSERVER

```
begin
...
/* socket fd is used for accepting service requests */
listen(fd, ...);
while (1) {
    /* Request served through socket newfd */
    newfd = accept(fd, ...);
    if (fork() == 0) { /* child */
        close(fd);
        /* serve request */
        ...
        exit(0);
    }
    else /* parent */
        close(newfd);
}
...
end
```

Figure 5: Concurrent server using sockets.

process. Second, the role of the network receive interrupt handler in our system is minimal, i.e., only to demultiplex packets to their destination processes. Receive side protocol processing is done in the context of the receive thread associated with a cached endpoint. Section 5 discusses ARC scheduling for protocol threads in Migrating Sockets.

4.2 Optimization for concurrent server model

In client/server programming, there are two principal programming models. They are the iterative server model and the concurrent server model. In the latter model, the server's role is only to listen for service requests from remote hosts. Once a request has been received, the server forks a child process to handle it, and itself goes back to listening for more requests. Because it allows new service requests to be accepted while previous ones are still being served, many programs, including the Internet Superserver (*inetd*), use the concurrent server model.

A typical program template for concurrent servers using sockets is shown in Figure 5. The program uses a socket *fd* to listen for incoming service requests. When a request arrives from a remote endpoint, it is accepted and a network connection is established. The local endpoint of the new connection is accessible through *newfd*. The server then forks a child process to serve the request and itself closes *newfd*.

In the caching mechanism described so far, *newfd* will be cached to the server on being returned by *accept()*. Immediately afterwards, however, the server does a *fork()* to serve the request in a child process, forcing *newfd* to be flushed in Migrating Sockets, because it is now shared between the server and the child. In fact, though, the server no longer needs access to *newfd*. When it closes *newfd*, the socket becomes cached to the child process.

Notice that although the final objective is to cache the accepted socket in the child process, two cache and one flush operations are involved, of which one cache and one flush would be unnecessary. To solve the problem, Migrating Sockets supports a new socket option called *SO_CONCURRENT_SERVER*. When the option is set

for a socket, say *S*, in the listen state, it serves as a hint that sockets accepted through *S* should ultimately be cached to a child process forked by the process, say *P*, doing the listen. Hence, when a socket is accepted through *S*, it is merely marked *cacheable* instead of being cached to *P*. If later, *P* does access *S* for send/receive, *S* will be cached (this is known as *delayed caching*). If, however, *P* does a fork before it accesses *S*, *S* will be cached to the forked process, say *Q*, as part of copying file descriptors from *P* to *Q*. Moreover, *S* is marked “recall-on-access” in the network server, meaning that if *P* later accesses *S*, the network server will recall *S* from *Q*, in effect causing *S* to be flushed. However, it is more likely that *P* will soon close *S*, and the recall-on-access status of *S* can be cleared.

4.3 Routing information management

The Internet protocol suite owes much of its flexibility and robustness in a heterogeneous and dynamic networking environment to protocols at (or below) the network layer. These protocols include, among others, IP, ICMP, IGMP, EGP and ARP. Together, they allow network routes to be dynamically discovered and reconfigured, and network connectivity is not lost even as network interfaces and routers come up or go down. In our system, we refer to the data structures that keep track of various kinds of routing information collectively as the *routing table*.

Dynamic routing interacts with caching of network endpoints. The reason is that even after a network connection has been established between a sender and receiver, the next “hop” (represented by an IP address) to which the sender must send its packets in order to reach the receiver may still change over time. This can happen, for example, by way of an ICMP redirect message, and the routing table must be updated to reflect the change. Moreover, the network interface to which some IP address has been assigned may be replaced by another interface, and ARP must update its translation of the IP address (in the routing table) to the link level address of the new interface.

In our design, we consider routing table management a global system function. As such, the network server is the only process responsible for its management. This has two advantages. First, routing table management functions do not have to be duplicated in the address space of every application process. Second, application processes do not need to be interrupted by (and process) routing messages that do not affect them. However, since application processes need read access to the routing table even for common case send and receive, such access must be as inexpensive as if the routing table were local to each process.

To satisfy the efficiency requirement, the network server creates a shared memory region, and allocates routing table entries exclusively from that region. Moreover, data pointers in the routing table must retain their intended meaning (without translation) irrespective of the process accessing them. This requires the network server and each application process to map² the shared memory region at a “well-known” virtual address. This virtual address can be returned by the network server to a client process at client registration time.

Our shared memory solution allows application processes to freely read the system routing table. We believe, however, that this does not represent a security problem in most cases. For example, users on a Unix system are often permitted to use the *netstat(1)* command to return the same kind of information.

²This region is mapped read-only by application processes.

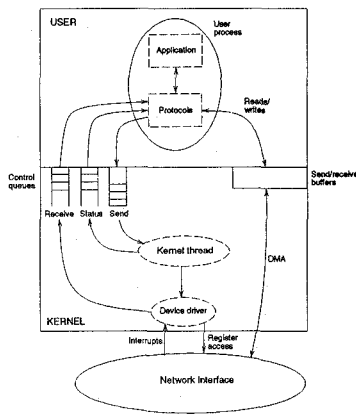


Figure 6: OS architecture for multimedia networking.

4.4 Protocol/kernel interface

For efficiency, Migrating Sockets runs on top of an OS architecture (Figure 6) we have previously prototyped for supporting continuous media (CM) applications [14].

Send/receive buffers shown in Figure 6 are allocated using the `IOBuffer::IOBuffer()` method in Table 1. The method creates a *network buffer* region for direct send/receive to/from the network (i.e. no intermediate data copies are required). If the network interface for send/receive uses DMA, the allocated network buffers will be automatically backed by required DMA resources. Moreover, buffers can be pinned in physical memory for predictable performance. A memory allocator associated with a buffer region supports flexible memory allocation/deallocation similar to `malloc(3C)` and `free(3C)` in the standard C library.

To send a packet, a user process appends control information in the form of a *send request* to a send control queue, managed through the `SendControlQueue` object in Table 1. The send request is then handled by the kernel through either a system call or a kernel thread introduced below. Notice that for some systems, even after kernel code processing of a send request has completed, the packet to send may only have been queued to a network interface, instead of really sent. Therefore, the buffer for sending cannot be reused or freed until the network interface has completed its part of the send and updated the status of the send request. That is why a method such as `reapall()` in Table 1 is necessary.

On the receive side, driver code for a network interface informs user processes of data to read by appending receive notifications to a receive control queue, managed through the `RecvControlQueue` object. For efficient control transfer, send/receive control queues are shared between user processes and the kernel.

A lightweight kernel thread provides shared access to a reserved-rate network connection in future integrated services networks. The kernel thread is periodically scheduled and implements a rate-based packet scheduling algorithm such that multiple user processes can send packets to the network connection with guaranteed data rates. For this purpose, a process creates a `MultiplexGroup` object using the `MultiplexGroup(int, int, int, void *, int)` method shown in Table 1. The method causes a *multiplex group*, a data structure used by the kernel thread for rate-based packet scheduling, to be created within the kernel. The `alg` parameter specifies the packet scheduling algorithm to use for the multiplex group. Currently, the `KT_RC` algorithm in [14] is supported. Param-

eters of the scheduling algorithm can be passed with the `params` pointer. For example, the `KT_RC` algorithm takes the scheduling period (in μs) of the kernel thread as a parameter. Once a multiplex group (identified by a key which is unique within an end system) has been created by a process, other processes can gain access to the group using the `MultiplexGroup(int)` method. Processes having access to a multiplex group can use the `join()` method to add traffic flows to the group with specified parameters. A rate parameter (in kbps), for example, is needed for a flow using the `KT_RC` algorithm. Flows can be deleted from a multiplex group using the `leave()` method.

Lastly, notice that, although the interface presented in Table 1 is designed to be general and device independent, some of the operations accessed through the interface are device dependent. For example, the `SendControlQueue::send()` method calls a device specific send function within the kernel. The `IOBuffer::IOBuffer()` method creates device dependent DMA resources backing allocated network buffers.

5 ARC Scheduling of Protocol Threads

Application and protocol threads in Migrating Sockets can specify their CPU requirements using the rate-based reservation model of ARC scheduling [15]. The rate-based model has two parameters: (1) rate, r , ($0 < r \leq 1$), and (2) period, p , in μs . Informally, the rate specifies a guaranteed fraction of CPU time that a thread with the reservation will be allocated over time intervals determined by p .

Progress guarantees of rate-based reservations are provided by instances of a family of ARC schedulers. ARC schedulers have the following properties: (1) reserved rate can be negotiated, (2) QoS guarantees are conditional upon thread behavior, and (3) firewall protection between threads is provided. The first property is provided by a *monitoring module* and a *rate-adaptation interface* as discussed in [15]. The second and third properties are provided by using an on-line CPU scheduling algorithm with the *firewall property*, such as the RC scheduler in [15]. Subject to the admission control criterion that the aggregate reserved rate in a system should not exceed one, RC provides the following progress guarantee: a “punctual” thread with rate r and period p is guaranteed at least $k r p$ CPU time over time interval $k p$, for $k = 1, 2, \dots$. Because ARC schedulers offer firewall protection between “well-behaved” and “greedy” threads, they are appropriate for integrated scheduling of continuous media and other applications found in a general purpose workstation.

In our current system, ARC scheduling has been implemented for multiplexing kernel threads onto the CPU, but not for a separate user level thread scheduler. In Solaris, however, user threads can be *bound* one-one to kernel threads. This allows ARC scheduling to work for all the protocol threads in Migrating Sockets.

6 Active Demultiplexing

An important task of protocol processing is to demultiplex incoming packets to their network endpoints. Traditionally, the receive side of a kernel level transport protocol looks for matches by searching a list of protocol control blocks known to the system. Recent user level protocol implementations have relied on packet filters installed with the driver of a network interface to accept or reject packets. While highly flexible, these methods involve searching. Although techniques such as hashing and one behind cache can significantly reduce the search time on the average, the actual search time may still be highly variable when the number of network endpoints in a

| Class | Method | Synopsis |
|------------------|---|--|
| IOBuffer | IOBuffer(int fd, caddr_t addr, int size); | Create a network buffer region of <u>size</u> bytes and map the region at user address <u>addr</u> . <u>fd</u> is a file descriptor for the network interface for which the buffer region is being allocated. |
| | void *malloc(int size); | Allocate a properly aligned buffer of <u>size</u> bytes from network buffer region. |
| | void free(void *buf); | Free to network buffer region buffer <u>buf</u> previously allocated by malloc(). |
| SendControlQueue | SendControlQueue(int fd, caddr_t addr, int size); | Allocate a send control queue of <u>size</u> send notifications for the network interface identified in <u>fd</u> . The send control queue is to be mapped at user address <u>addr</u> . |
| | void attach(IOBuffer *iob); | Attach network buffer region <u>iob</u> to send control queue. |
| | int send(void *buf, int len); | Send <u>len</u> bytes starting at user address <u>buf</u> by appending a send notification to send control queue. |
| | struct status *reap(int block); | Return pointer to status of last newly completed send. |
| | void reapall(); | Free buffers in all newly completed sends by calling free() method of attached IOBuffer. |
| RecvControlQueue | RecvControlQueue(int fd, caddr_t addr, int size); | Allocate a receive control queue of <u>size</u> receive notifications and map it at user address <u>addr</u> . <u>fd</u> is a file descriptor for the network interface for which control queue is allocated. |
| | int recv(caddr_t *buf, int *len, int block); | Receive <u>len</u> bytes of data in network buffer whose pointer is returned in <u>buf</u> . If no data are available for receive, block calling thread iff <u>block</u> is set. |
| | void unmap(caddr_t buf); | Give back buffer <u>buf</u> previously returned by recv() to kernel. |
| MultiplexGroup | MultiplexGroup(int key, int bw, int alg, void *params, int perm); | Create multiplex group with key <u>key</u> for a reserved rate network connection, with bandwidth <u>bw</u> Mbps. Multiplexing algorithm is specified by <u>alg</u> , with algorithm parameters pointed to by <u>params</u> . Permission for other processes to access multiplex group is specified in <u>perm</u> . |
| | MultiplexGroup(int key); | Get a multiplex group previously created with key <u>key</u> . |
| | int setoption(int optname, void *optval); | Set option for multiplex group. Option name is <u>optname</u> and value is pointed to by <u>optval</u> . |
| | int join(void *params); | Let a flow join multiplex group with parameters pointed to by <u>params</u> . Return id for the flow in multiplex group. |
| | void leave(int id); | Delete flow with id <u>id</u> from multiplex group. |

Table 1: Kernel interface for protocol code in Migrating Sockets.

system is large.

Since predictable performance is an important goal of our architecture, our system supports a constant overhead packet demultiplexing mechanism known as *active demultiplexing*. The basic idea is that, under certain conditions, an OS handle identifying a receive process can be included in packets destined for that process. An end system can then make use of the OS handle to deliver packets directly to the receive process, without any searching.

Currently, active demultiplexing exploits the notion of exclusive packet receivers introduced in section 3. In situations where active demultiplexing cannot be applied or is not preferred, our system provides packet filters.

6.1 Mechanism

We have implemented active demultiplexing in the context of TCP/IP. The mechanism can be transparently enabled when both sender and receiver hosts in a TCP connection support it. Figure 7 illustrates the mechanism, triggered when a socket, say *S*, that is an exclusive packet receiver becomes cached to a user process. An OS handle for the user process to receive from *S* becomes known as part of the process of caching. The newly cached socket *S* advertises the OS handle when it next sends a packet to the remote endpoint, say *R*, of the network connection. This happens when

either *S* has data to send to *R*, or when *S* acknowledges packets received from *R*. The advertisement is carried in a new TCP option called TCPOPT_DEMUX_ADVERTISE.

On processing TCPOPT_DEMUX_ADVERTISE, *R* learns about *S*'s OS handle. *R* then caches the handle and can later *enclose* it in subsequent packets it sends to *S*, by way of a new IP option IPOPT_DEMUX_ENCLOSED. The enclosed OS handle enables active demultiplexing by *S*'s end system. To allow a link level device driver to easily locate IPOPT_DEMUX_ENCLOSED, the option is always inserted as a first IP option. In addition, we use a currently unused bit in the service field of IP header to indicate to the device driver whether an OS handle has been enclosed.

An OS handle should be revoked when a cached socket becomes no longer exclusive, or when an exclusive socket becomes flushed. Handle revocation is achieved using a new TCP option TCPOPT_DEMUX_REVOKE. On receiving a handle revocation from its peer, a network endpoint stops including IPOPT_DEMUX_ENCLOSED in packets sent to the peer. Figure 8 shows the formats of the various options used for active demultiplexing. Notice that a sequence number is included in handle advertisements and revocations. We prescribe that the two operations be applied only in increasing sequence number order, thereby preventing an earlier operation from overwriting a more recent one.

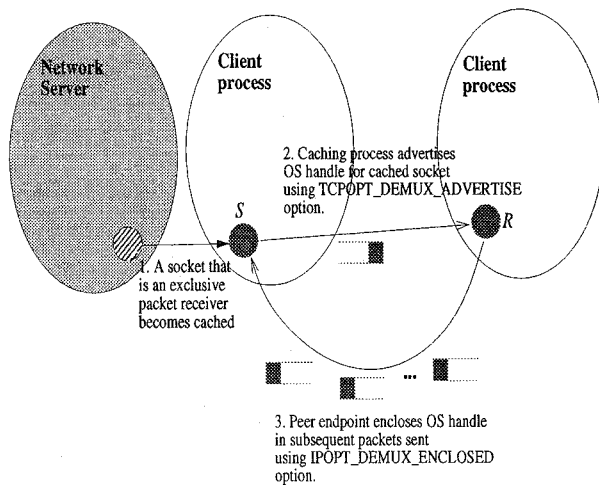


Figure 7: Mechanism of active demultiplexing for TCP/IP.

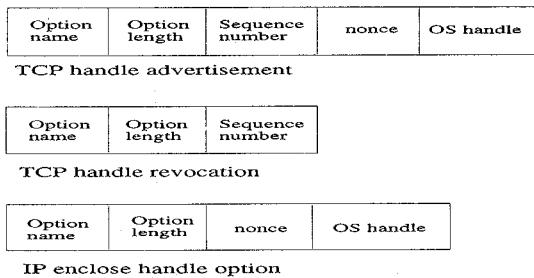


Figure 8: Formats of active demultiplexing options.

The complete state transition diagram for handle advertisement and revocation using TCP/IP is shown in Figure 9. If (and only if) an endpoint is in the “advertising” state, the advertisement option will be used for TCP segments sent by the endpoint. Similarly, if (and only if) an endpoint is in the “revoking” state, the revocation option will be used for TCP segments sent by the endpoint.

6.2 Security considerations

Taking an OS handle in an incoming packet and using it to deliver the packet directly to a receive process is a “powerful” mechanism. Without proper precautions, the mechanism can raise some serious security concerns. First, a malicious process can advertise an indiscriminate OS handle (i.e. one that is not associated with an exclusive packet receiver) and “hoard” packets that should also be delivered to other processes in the system. Second, an OS handle can become obsolete, such as when a cached endpoint becomes flushed. Third, a faulty process can enclose a wrong or fabricated OS handle in sending packets.

To guard against such security problems, we generate a *nonce* when a network endpoint is installed and associate it with the endpoint. The nonce is used in conjunction with the OS handle for active demultiplexing. Nonces have the properties that each newly generated nonce has a fresh value, and it is difficult to guess the value of a nonce that is not explicitly passed. Before accepting an OS handle, receive side demultiplexing code performs security checks as shown in algorithm ACTIVE (Figure 10). Notice that the check at line 2 guards against the first security threat in the preceding paragraph, while the checks at lines 1 and 3 guard against the second and third security threats. We fall back on a conventional packet

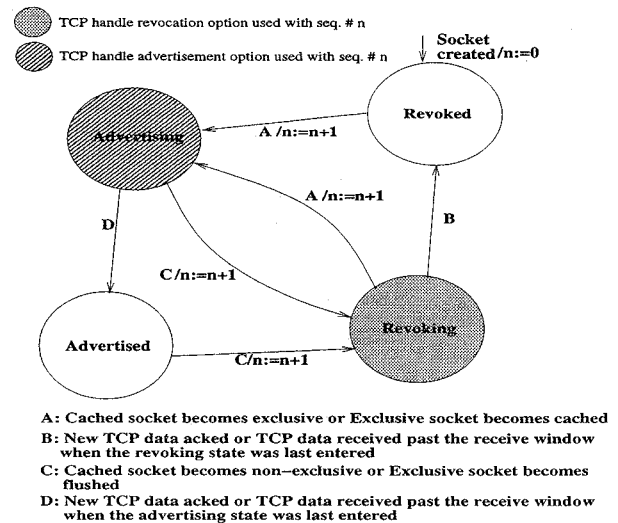


Figure 9: State transition diagram for handle advertisement and revocation in TCP/IP.

Algorithm ACTIVE

begin

- Remember enough execution state to transfer to fallback on data fault;
 - if (handle is not for an exclusive packet receiver)
goto fallback;
 - if (nonce in handle does not match nonce for receive endpoint)
goto fallback;
 - Deliver packet directly to process identified in handle;
 - return;
- fallback:
- Match packet against installed packet filters in system;

end ;

Figure 10: Specification of active demultiplexing algorithm.

filter mechanism if an enclosed OS handle is found unacceptable.

Finally, we note that, in practice, security can be better enforced if an OS handle is inserted by kernel level code (i.e. by a link level device driver supporting active demultiplexing). In this case, security safeguard is mainly to limit the cycle time of nonces, and it is sufficient to generate a new nonce by incrementing it cyclically.

7 Experimental Results

We have an implementation of Migrating Sockets on Solaris 2.5. We are currently running it on SPARC Ultra-1/140, SPARC 10/30 and SPARC 20/50 machines interconnected by a 10 Mbps Ethernet network. Kernel level support as discussed in section 4.4 is running on the Ultra-1 and SPARC 10, but not currently on the SPARC 20. The ARC scheduling framework presented in section 5 is being ported from Solaris 2.3 to Solaris 2.5. Major Internet protocols including TCP, UDP, and IP (with IP multicast) have been ported from 4.4 BSD to Migrating Sockets. For UDP, we have incorporated the following optimization techniques proposed in [11]: (1) Integrated

| Code | TCP send | | UDP send | |
|---------------|------------|--------|------------|--------|
| | 1400 bytes | 1 byte | 1400 bytes | 1 byte |
| Socket send | 90(74) | 72(62) | 16(9) | 15(9) |
| TCP/UDP | 144(73) | 25(4) | 101(62) | 28(7) |
| IP | 10(8) | 10(2) | 13(3) | 13(3) |
| Link layer | 9(2) | 9(2) | 15(3) | 14(3) |
| Device driver | 20 | 19 | 20 | 20 |
| Total | 273 | 135 | 165 | 90 |

Table 2: Breakdown of Migrating Sockets TCP/UDP send path latency (μ s) for Ultra-1.

checksumming and copying of data from application buffers to network buffers, (2) replacement of general purpose socket send code with more efficient UDP specific code, and (3) delete of pseudo-connect in UDP send. Moreover, since Internet checksumming is heavily used, we replaced the 4.4 BSD checksum routine with a more efficient routine optimized for the SPARC Ultra-1 architecture. Apart from optimized checksumming, the TCP code is largely ported as is from 4.4 BSD.

7.1 TCP/UDP performance

We measured the performance of the current implementation of TCP/IP and UDP/IP protocol stacks accessed through a cached socket in Migrating Sockets. Table 2 gives a breakdown of the various (average) component costs (the numbers in brackets are corresponding standard deviations) due to host software on the send path of a packet. Each packet carried either 1400 bytes or one byte of user data. The relatively high costs of socket and TCP send code for 1400 bytes were mainly due to data copy and TCP checksumming. For UDP, socket send code had minimal cost because it called a UDP specific send function very early on. The UDP function performed integrated checksum and copy of data from application to network buffers. Performance benefits of replacing “baroque” socket send code with protocol specific code can also be seen by comparing the costs of socket send for one byte of TCP and UDP data, respectively.

The row labeled “Link layer” in Table 2 gives the cost of link layer processing such as ARP address translation by Migrating Sockets. Also notice from the last row of the table that the kernel level overhead, i.e., processing by the network device driver, was largely insensitive to the packet size. This is because protocol code allocated network buffers for sending and data did not have to be copied from user to kernel space.

Table 3 gives a breakdown of the various (average) component costs (the numbers in brackets are the corresponding standard deviations) due to host software on a TCP/IP or UDP/IP packet receive path. As for the send path, a packet carried either 1400 bytes or one byte of user data. The “Kernel interrupt” number is the total time spent in the interrupt handler of the receive network device driver, which includes the cost of matching a packet to a BSD packet filter (which took roughly 10 μ s). A matched packet woke up a receive thread of Migrating Sockets and caused the thread to be scheduled. The cost of context switching to the receive thread is given by the “Switch to recv thread” number.

Protocol processing of a received packet by Migrating Sockets starts at the link layer. The “Link layer” number includes costs such as inspecting the type field of the Ethernet header and dispatching the packet to IP. TCP or UDP processing of a packet took longer for 1400 bytes than for one byte of user data, due to checksumming of a longer packet. After the transport layer, the receive thread ap-

| Code | TCP receive | | UDP receive | |
|-----------------------|-------------|--------|-------------|--------|
| | 1400 B | 1 B | 1400 B | 1 B |
| Return from read | 71(28) | 52(24) | 45(24) | 30(11) |
| Switch to read thread | 79(16) | 92(51) | 56(18) | 56(16) |
| TCP/UDP | 90 | 56 | 89 | 57 |
| IP | 15 | 14 | 9(3) | 10(3) |
| Link layer | 8(7) | 7(5) | 6(3) | 7(3) |
| Switch to recv thread | 34 | 33 | 38 | 32 |
| Kernel interrupt | 33 | 32 | 33 | 31 |
| Total | 330 | 286 | 276 | 223 |

Table 3: Breakdown of Migrating Sockets TCP/UDP receive path latency (μ s) for Ultra-1.

ended the packet to a socket receive buffer and woke up an upper half application thread that was blocked reading from the socket. The “Switch to read thread” number gives the overhead of context switching to the application thread. Finally, the “Return from read” number includes the cost of copying any user data from a socket receive buffer to an application buffer. Hence, it was higher for a larger packet size.

We note that although high performance is not the main concern in our work, our TCP/UDP latency numbers do seem to compare very well with those reported in the literature (e.g. [8]).³

7.2 Optimized checksumming

This set of experiments quantifies the performance benefits of using an Internet checksum routine that is optimized for a specific computer architecture, namely the SPARC Ultra-1. Figure 11 compares the performance between the original 4.4 BSD checksum routine and an optimized checksum routine for various data sizes. By exploiting knowledge of the most efficient data size and alignment for memory access, the optimized routine achieves significant improvement over the original one.

We next investigate the performance benefits of integrated checksum and copy. In the experiments, the destination address of a copy was double-word aligned. For the separate checksum and copy approach, the copy loop was done immediately after the checksum loop. Figure 12 shows the results when the test program had a small memory footprint (less than 8 M-Bytes). For the integrated approach, the performance was very slightly better if the source address was also double-word aligned than if it was word or byte-aligned. Also, the integrated approach showed minimal improvement over the separate approach because of cache effects. Figure 13 shows the results when the test program had a large memory footprint (41 M-Bytes). In this case, the integrated approach was significantly more efficient than the separate approach.

7.3 Active demultiplexing

Table 4 shows the per packet active demultiplexing overhead in our current system. For comparison, we note that it took about 10 μ s to match a packet with a BSD packet filter on the Ultra-1. Notice that the kernel level overhead for executing Algorithm ACTIVE (Figure 10) is quite small. The dominant cost of inserting IPOPT_DEMUX_ENCLOSED option (of size 8 bytes) in an IP packet is due to the fact that our current system treats inserting any IP option as an “exceptional” case. Hence, every time IPOPT_DEMUX_ENCLOSED is to be inserted in a packet, the TCP/IP headers of the packet must be moved to make room

³Our numbers are smaller. However, a direct comparison is unfair because of the use of different hardware platforms.

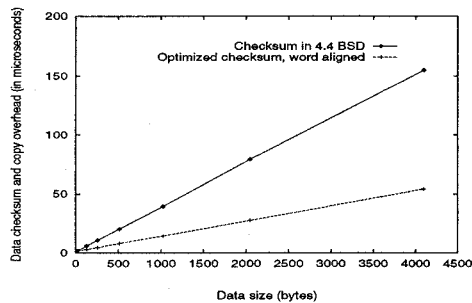


Figure 11: Checksum overhead (in μs) versus data size (in bytes) for Ultra-1.

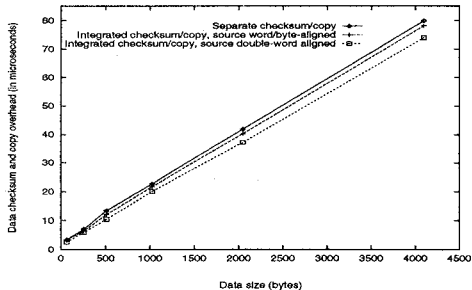


Figure 12: Data checksum and copy overhead (in μs) versus data size (in bytes) for small memory footprint (Ultra-1).

for the option. A straightforward optimization is to include IPOPT_DEMUX_ENCLOSED as part of the header template for packets to be sent from relevant sockets.

| Machine | Insert IP option | Data fault protection | Other kernel overhead | Total |
|----------|------------------|-----------------------|-----------------------|-------|
| SPARC 10 | 5.9 | 0.702 | 0.256 | 6.858 |
| Ultra-1 | 2.5 | 0.225 | 0.056 | 2.781 |

Table 4: Breakdown of per packet active demultiplexing overhead (in μs).

8 Conclusion

We presented Migrating Sockets as a framework for user level protocol implementation, and discussed its relations to other system components in an end system architecture for networking with QoS guarantees.

Acknowledgment

The authors wish to thank Long Ma for helping with performance measurements of the system.

REFERENCES

1. S. Bradner and A. Mankin. The recommendation for the IP next generation protocol. Internet RFC 1752, January 1995.
2. Torsten Braun and Christophe Diot. Protocol implementation using integrated layer processing. In *Proc. ACM SIGCOMM '95*, Boston, MA, August 1995.
3. Chris Dalton, Greg Watson, David Banks, Costas Calamvokis, Aled Edwards, and John Lumley. Afterburner. *IEEE Network*, 7(4):36–43, July 1993.

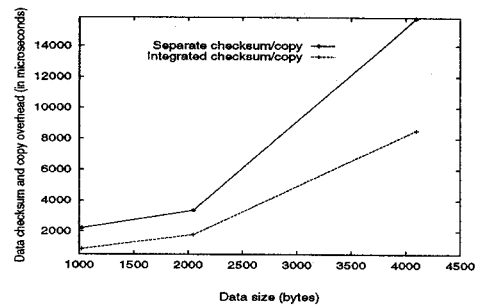


Figure 13: Data checksum and copy overhead (in μs) versus data size (in bytes) for large memory footprint (Ultra-1).

4. D.R. Engler and M.F. Kaashoek. DPF: fast, flexible message demultiplexing using dynamic code generation. In *Proc. ACM SIGCOMM '96*, Stanford, CA, August 1996.
5. ATM Forum. ATM traffic management specification, version 4.0, 1995.
6. R. Gopalakrishnan and G.M. Parulkar. Real-time upcalls: A mechanism to provide real-time processing guarantees. Technical report, Washington University, 1995.
7. N.C. Hutchinson, S. Mishra, L.L. Peterson, and V.T. Thomas. Tools for implementing network protocols. *Software - Practice and Experience*, 1989.
8. Chris Maeda and Brian N. Bershad. Protocol service decomposition for high-performance networking. In *Proc. 14th SOSP*, pages 244–255, December 1993.
9. S. McCanne and Van Jacobson. The BSD packet filter: A new architecture for user-level packet capture. In *USENIX Technical Conference Proceedings*, pages 259–269, San Diego, CA, Winter 1993.
10. Clifford W. Mercer, Jim Zelenka, and Ragunathan Rajkumar. On predictable operating system protocol processing. Technical Report CMU-CS-94-165, Carnegie Mellon University, Pittsburgh, PA, May 1994.
11. C. Partridge and S. Pink. A faster UDP. *IEEE/ACM Transactions on Networking*, 1(4):429–440, August 1993.
12. C.A. Thekkath, T.D. Nguyen, E. Moy, and E.D. Lazowska. Implementing network protocols at user level. *IEEE/ACM Trans. Networking*, 1(5):554–565, October 1993.
13. Thorsten von Eicken, Anindya Basu, Vineet Buch, and Werner Vogels. U-Net: A user-level network interface for parallel and distributed computing. In *Proc. 15th SOSP*, November 1995.
14. David K.Y. Yau and Simon S. Lam. An architecture towards efficient OS support for distributed multimedia. In *Proc. IS&T/SPIE Multimedia Computing and Networking*, pages 424–435, San Jose, CA, January 1996.
15. David K.Y. Yau and Simon S. Lam. Adaptive rate-controlled scheduling for multimedia applications. *IEEE/ACM Transactions on Networking*, August 1997.
16. Lixia Zhang, Stephen Deering, Deborah Estrin, Scott Shenker, and Daniel Zappala. RSVP: A new resource ReSerVation Protocol. *IEEE Network*, pages 8–18, September 1993.