# WRL
# Research Report 95/8

# Eliminating
# Receive Livelock
# in an
# Interrupt-driven Kernel

*Jeffrey C. Mogul*
*K. K. Ramakrishnan*

The Western Research Laboratory (WRL) is a computer systems research group that was founded by Digital Equipment Corporation in 1982. Our focus is computer science research relevant to the design and application of high performance scientific computers. We test our ideas by designing, building, and using real systems. The systems we build are research prototypes; they are not intended to become products.

There are two other research laboratories located in Palo Alto, the Network Systems Lab (NSL) and the Systems Research Center (SRC).  Another Digital research group is located in Cambridge, Massachusetts (CRL).

Our research is directed towards mainstream high-performance computer systems. Our prototypes are intended to foreshadow the future computing environments used by many Digital customers. The long-term goal of WRL is to aid and accelerate the development of high-performance uni- and multi-processors. The research projects within WRL will address various aspects of high-performance computing.

We believe that significant advances in computer systems do not come from any single technological advance. Technologies, both hardware and software, do not all advance at the same pace. System design is the art of composing systems which use each level of technology in an appropriate balance. A major advance in overall system performance will require reexamination of all aspects of the system.

We do work in the design, fabrication and packaging of hardware; language processing and scaling issues in system software design; and the exploration of new applications areas that are opening up with the advent of higher performance systems. Researchers at WRL cooperate closely and move freely among the various levels of system design. This allows us to explore a wide range of tradeoffs to meet system goals.

We publish the results of our work in a variety of journals, conferences, research reports, and technical notes.  This document is a research report. Research reports are normally accounts of completed research and may include material from earlier technical notes.  We use technical notes for rapid distribution of technical material; usually this represents research in progress.

Research reports and technical notes may be ordered from us.  You may mail your order to:

Technical Report Distribution
DEC Western Research Laboratory, WRL-2
250 University Avenue
Palo Alto, California 94301   USA

Reports and technical notes may also be ordered by electronic mail. Use one of the following addresses:

Digital E-net:          `JOVE::WRL-TECHREPORTS`

Internet:          `WRL-Techreports@decwrl.pa.dec.com`

UUCP:          `decpa!wrl-techreports`

To obtain more details on ordering by electronic mail, send a message to one of these addresses with the word ''`help`'' in the Subject line; you will receive detailed instructions.

Reports and technical notes may also be accessed via the World Wide Web: `http://www.research.digital.com/wrl/home.html`.

# Eliminating Receive Livelock
# in an Interrupt-driven Kernel

## Jeffrey C. Mogul

Digital Equipment Corporation Western Research Laboratory
mogul@wrl.dec.com

## K. K. Ramakrishnan

AT&T Bell Laboratories
600 Mountain Avenue, Murray Hill, New Jersey 07974
kkrama@research.att.com
(work done while at Digital Equipment Corporation)

**December, 1995**

## Abstract

Most operating systems use interface interrupts to schedule network tasks. Interrupt-driven systems can provide low overhead and good latency at low offered load, but degrade significantly at higher arrival rates unless care is taken to prevent several pathologies. These are various forms of ***receive livelock***, in which the system spends all its time processing interrupts, to the exclusion of other necessary tasks. Under extreme conditions, no packets are delivered to the user application or the output of the system.

To avoid livelock and related problems, an operating system must schedule network interrupt handling as carefully as it schedules process execution. We modified an interrupt-driven networking implementation to do so; this eliminates receive livelock without degrading other aspects of system performance. We present measurements demonstrating the success of our approach.

This report is an expanded version of a paper in the *Proceedings of the 1996 USENIX Technical Conference*.

# Table of Contents

# List of Figures

# List of Tables

# 1. Introduction

Most operating systems use interrupts to internally schedule the performance of tasks related to I/O events, and particularly the invocation of network protocol software. Interrupts are useful because they allow the CPU to spend most of its time doing useful processing, yet respond quickly to events without constantly having to poll for event arrivals.

Polling is expensive, especially when I/O events are relatively rare, as is the case with disks, which seldom interrupt more than a few hundred times per second. Polling can also increase the latency of response to an event. Modern systems can respond to an interrupt in a few tens of microseconds; to achieve the same latency using polling, the system would have to poll tens of thousands of times per second, which would create excessive overhead. For a general-purpose system, an interrupt-driven design works best.

Most extant operating systems were designed to handle I/O devices that interrupt every few milliseconds. Disks tended to issue events on the order of once per revolution; first-generation LAN environments tend to generate a few hundred packets per second for any single end-system. Although people understood the need to reduce the cost of taking an interrupt, in general this cost was low enough that any normal system would spend only a fraction of its CPU time handling interrupts.

The world has changed. Operating systems typically use the same interrupt mechanisms to control both network processing and traditional I/O devices, yet many new applications can generate packets several orders of magnitude more often than a disk can generate seeks. Multimedia and other real-time applications will become widespread. Client-server applications, such as NFS, running on fast clients and servers can generate heavy RPC loads. Multicast and broadcast protocols subject innocent-bystander hosts to loads that do not interest them at all. As a result, network implementations must now deal with significantly higher event rates.

Many multi-media and client-server applications share another unpleasant property: unlike traditional network applications (Telnet, FTP, electronic mail), they are not flow-controlled. Some multi-media applications want constant-rate, low-latency service; RPC-based client-server applications often use datagram-style transports, instead of reliable, flow-controlled protocols. Note that whereas I/O devices such as disks generate interrupts only as a result of requests from the operating system, and so are inherently flow-controlled, network interfaces generate unsolicited receive interrupts.

The shift to higher event rates and non-flow-controlled protocols can subject a host to congestive collapse: once the event rate saturates the system, without a negative feedback loop to control the sources, there is no way to gracefully shed load. If the host runs at full throughput under these conditions, and gives fair service to all sources, this at least preserves the possibility of stability. But if throughput decreases as the offered load increases, the overall system becomes unstable.

Interrupt-driven systems tend to perform badly under overload. Tasks performed at interrupt level, by definition, have absolute priority over all other tasks. If the event rate is high enough to cause the system to spend all of its time responding to interrupts, then nothing else will happen, and the system throughput will drop to zero. We call this condition *receive livelock*: the system is not deadlocked, but it makes no progress on any of its tasks.

Any purely interrupt-driven system using fixed interrupt priorities will suffer from receive livelock under input overload conditions. Once the input rate exceeds the reciprocal of the CPU cost of processing one input event, any task scheduled at a lower priority will not get a chance to run.

Yet we do not want to lightly discard the obvious benefits of an interrupt-driven design. Instead, we should integrate control of the network interrupt handling sub-system into the operating system's scheduling mechanisms and policies. In this paper, we present a number of simple modifications to the purely interrupt-driven model, and show that they guarantee throughput and improve latency under overload, while preserving the desirable qualities of an interrupt-driven system under light load.

## 2. Motivating applications

We were led to our investigations by a number of specific applications that can suffer from livelock. Such applications could be built on dedicated single-purpose systems, but are often built using a general-purpose system such as UNIX®, and we wanted to find a general solution to the livelock problem. The applications include:

- *Host-based routing*: Although inter-network routing is traditionally done using special-purpose (usually non-interrupt-driven) router systems, routing is often done using more conventional hosts. Virtually all Internet ''firewall'' products use UNIX or Windows NT™ systems for routing [12, 19]. Much experimentation with new routing algorithms is done on UNIX [6], especially for IP multicasting.

- *Passive network monitoring*: network managers, developers, and researchers commonly use UNIX systems, with their network interfaces in ''promiscuous mode,'' to monitor traffic on a LAN for debugging or statistics gathering [13].

- *Network file service*: servers for protocols such as NFS are commonly built from UNIX systems.

These applications (and others like them, such as Web servers) are all potentially exposed to heavy, non-flow-controlled loads.

We have encountered livelock in all three of these applications, have solved or mitigated the problem, and have shipped the solutions to customers. For example, an early implementation of this work was successfully deployed in the routers used for the NASDAQ financial network.

The rest of this paper concentrates on host-based routing and (to a lesser extent) network monitoring, since this simplifies the context of the problem and allows easy performance measurement.

## 3. Requirements for scheduling network tasks

Performance problems generally arise when a system is subjected to transient or long-term input overload. Ideally, the communication subsystem could handle the worst-case input load without saturating, but cost considerations often prevent us from building such powerful systems. Systems are usually sized to support a specified design-center load, and under overload the best we can ask for is controlled and graceful degradation.

When an end-system is involved in processing considerable network traffic, its performance depends critically on how its tasks are scheduled. The mechanisms and policies that schedule packet processing and other tasks should guarantee acceptable system *throughput*, reasonable *latency* and *jitter* (variance in delay), *fair* allocation of resources, and overall system *stability*, without imposing excessive overheads, especially when the system is overloaded.

We can define throughput as the rate at which the system delivers packets to their ultimate consumers. A consumer could be an application running on the receiving host, or the host could be acting as a router and forwarding packets to consumers on other hosts. We expect the throughput of a well-designed system to keep up with the offered load up to a point called the *Maximum Loss Free Receive Rate* (MLFRR), and at higher loads throughput should not drop below this rate.

Of course, useful throughput depends not just on successful reception of packets; the system must also transmit packets. Because packet reception and packet transmission often compete for the same resources, under input overload conditions the scheduling subsystem must ensure that packet transmission continues at an adequate rate.

Many applications, such as distributed systems and interactive multimedia, often depend more on low-latency, low-jitter communications than on high throughput. Even during overload, we want to avoid long queues, which increases latency, and bursty scheduling, which increases jitter.

When a host is overloaded with incoming network packets, it must also continue to process other tasks, so as to keep the system responsive to management and control requests, and to allow applications to make use of the arriving packets. The scheduling subsystem must fairly allocate CPU resources among packet reception, packet transmission, protocol processing, other I/O processing, system housekeeping, and application processing.

A host that behaves badly when overloaded can also harm other systems on the network. Livelock in a router, for example, may cause the loss of control messages, or delay their processing. This can lead other routers to incorrectly infer link failure, causing incorrect routing information to propagate over the entire wide-area network. Worse, loss or delay of control messages can lead to network instability, by causing positive feedback in the generation of control traffic [15].

## 4. Interrupt-driven scheduling and its consequences

Scheduling policies and mechanisms significantly affect the throughput and latency of a system under overload. In an interrupt-driven operating system, the interrupt subsystem must be viewed as a component of the scheduling system, since it has a major role in determining what code runs when. We have observed that interrupt-driven systems have trouble meeting the requirements discussed in section 3.

In this section, we first describe the characteristics of an interrupt-driven system, and then identify three kinds of problems caused by network input overload in interrupt-driven systems:

- *Receive livelocks* under overload: delivered throughput drops to zero while the input overload persists.

- Increased *latency* for packet delivery or forwarding: the system delays the delivery of one packet while it processes the interrupts for subsequent packets, possibly of a burst.

- *Starvation* of packet transmission: even if the CPU keeps up with the input load, strict priority assignments may prevent it from transmitting any packets.

## 4.1. Description of an interrupt-driven system

An interrupt-driven system performs badly under network input overload because of the way in which it prioritizes the tasks executed as the result of network input. We begin by describing a typical operating system's structure for processing and prioritizing network tasks. We use the 4.2BSD [9] model for our example, but we have observed that other operating systems, such as VMS™, DOS, and Windows NT, and even several Ethernet chips, have similar characteristics and hence similar problems.

When a packet arrives, the network interface signals this event by interrupting the CPU. Device interrupts normally have a fixed Interrupt Priority Level (IPL), and preempt all tasks running at a lower IPL; interrupts do not preempt tasks running at the same IPL. The interrupt causes entry into the associated network device driver, which does some initial processing of the packet. In 4.2BSD, only buffer management and data-link layer processing happens at ''device IPL.'' The device driver then places the packet on a queue, and generates a software interrupt to cause further processing of the packet. The software interrupt is taken at a lower IPL, and so this protocol processing can be preempted by subsequent interrupts. (We avoid lengthy periods at high IPL, to reduce latency for handling certain other events.)

The queues between steps executed at different IPLs provide some insulation against packet losses due to transient overloads, but typically they have fixed length limits. When a packet should be queued but the queue is full, the system must drop the packet. The selection of proper queue limits, and thus the allocation of buffering among layers in the system, is critical to good performance, but beyond the scope of this paper.

Note that the operating system's scheduler does not participate in any of this activity, and in fact is entirely ignorant of it.

As a consequence of this structure, a heavy load of incoming packets could generate a high rate of interrupts at device IPL. Dispatching an interrupt is a costly operation, so to avoid this overhead, the network device driver attempts to *batch* interrupts. That is, if packets arrive in a burst, the interrupt handler attempts to process as many packets as possible before returning from the interrupt. This amortizes the cost of processing an interrupt over several packets.

Even with batching, a system overloaded with input packets will spend most of its time in the code that runs at device IPL. That is, the design gives absolute priority to processing incoming packets. At the time that 4.2BSD was developed, in the early 1980s, the rationale for this was that network adapters had little buffer memory, and so if the system failed to move a received packet promptly into main memory, a subsequent packet might be lost. (This is still a problem with low-cost interfaces.) Thus, systems derived from 4.2BSD do minimal processing at device IPL, and give this processing priority over all other network tasks.

Modern network adapters can receive many back-to-back packets without host intervention, either through the use of copious buffering or highly autonomous DMA engines. This insulates the system from the network, and eliminates much of the rationale for giving absolute priority to the first few steps of processing a received packet.

## 4.2. Receive livelock

In an interrupt-driven system, receiver interrupts take priority over all other activity. If packets arrive too fast, the system will spend all of its time processing receiver interrupts. It will therefore have no resources left to support delivery of the arriving packets to applications (or, in the case of a router, to forwarding and transmitting these packets). The useful throughput of the system will drop to zero.

Following [16], we refer to this condition as *receive livelock*: a state of the system where no useful progress is being made, because some necessary resource is entirely consumed with processing receiver interrupts. When the input load drops sufficiently, the system leaves this state, and is again able to make forward progress. This is not a deadlock state, from which the system would not recover even when the input rate drops to zero.

A system could behave in one of three ways as the input load increases. In an ideal system, the delivered throughput always matches the offered load. In a realizable system, the delivered throughput keeps up with the offered load up to the *Maximum Loss Free Receive Rate* (MLFRR), and then is relatively constant after that. At loads above the MLFRR, the system is still making progress, but it is dropping some of the offered input; typically, packets are dropped at a queue between processing steps that occur at different priorities.

In a system prone to receive livelock, however, throughput decreases with increasing offered load, for input rates above the MLFRR. Receive livelock occurs at the point where the throughput falls to zero. A livelocked system wastes all of the effort it puts into partially processing received packets, since they are all discarded.

Receiver-interrupt batching complicates the situation slightly. By improving system efficiency under heavy load, batching can increase the MLFRR. Batching can shift the livelock point but cannot, by itself, prevent livelock.

In section 6.2, we present measurements showing how livelock occurs in a practical situation. Additional measurements, and a more detailed discussion of the problem, are given in [16].

## 4.3. Receive latency under overload

Although interrupt-driven designs are normally thought of as a way to reduce latency, they can actually increase the latency of packet delivery. If a burst of packets arrives too rapidly, the system will do link-level processing of the entire burst before doing any higher-layer processing of the first packet, because link-level processing is done at a higher priority. As a result, the first packet of the burst is not delivered to the user until link-level processing has been completed for all the packets in the burst. The latency to deliver the first packet in a burst is increased almost by the time it takes to receive the entire burst. If the burst is made up of several independent

NFS RPC requests, for example, this means that the server's disk sits idle when it could be doing useful work.

To demonstrate this effect, we performed experiments using ULTRIX™ Version 3.0 running on a DECstation 3100 (approximately 11.3 SPECmarks). ULTRIX, derived from 4.2BSD, closely follows the network design of that system. We used a logic analyzer to measure the time between the generation of an interrupt by the Ethernet device (an AMD 7990 LANCE chip), signalling the complete reception of a packet, and the packet's delivery to an application. We used the kernel's implementation of a simple data-link layer protocol, rather than IP/TCP or a similar protocol stack, but the steps performed by the kernel are substantially the same:

- link-level processing at device IPL, which includes copying the packet into kernel buffers

- further processing following a software interrupt, which includes locating the appropriate user process, and and queueing the packet for delivery to this process

- finally, awakening the user process, which then (in kernel mode) copies the received packet into its own buffer.

Figure 4-1 shows a time line for the completion of these processing stages, when receiving a burst of four minimum size packets from the Ethernet. The system starts to copy the first packet into a kernel buffer almost immediately after it arrives, but does not finish copying the third packet until about 1.33 msec later. Only after finishing this does it schedule a software interrupt to dispatch the packet to the user process, and all of the packets are dispatched before the user process is awakened. It is the use of preemptive interrupt priorities that prevents completion of processing for the first packet until substantial processing has been done on the entire burst.



**Figure 4-1:** How interrupt-driven scheduling causes excess latency under overload

We generated our bursts of Ethernet packets with an inter-packet spacing of 108 usec (this is not the minimum theoretical spacing, but we were limited by the packet generator we used). The latency to deliver the first packet to the user application depended on the size of a burst: 1.23 msec for a single-packet burst, 1.54 msec for a two-packet burst, 2.02 msec for a four-packet burst, and 5.03 msec for a 16-packet burst. Figure 4-2 summarizes this data, showing that the excess delay is nearly linear in the burst size.

Figure 4-2 also shows that the linear behavior persists for various larger packet sizes. The horizontal axis shows the number of packets in a burst; the vertical axis shows delivery latency

for the *first* packet in a burst. The data-link layer facility used in these trials supports only a limited amount of buffering, so in some trials not all of the packets in a burst were delivered to the application. The delivery queue could only hold 10 packets, so in the 16-packet trials only the first ten packets of a burst were actually received. These trials are indicated by open circles. Also, the queue could hold only 4 Kbytes of total data, so for the two trials using maximum-size packets, the system only delivered two packets from each burst; these trials are indicated by open squares.



*In trials marked with open squares or open circles, the kernel failed to deliver the entire burst to the user application.*

**Figure 4-2:** Receive latency as a function of burst size and packet size

We will present a more detailed analysis of receive latency in section 7.2, in the context of a somewhat different system.

## 4.4. Starvation of transmits under overload

In most systems, the packet transmission process consists of selecting packets from an output queue, handing them to the interface, waiting until the interface has sent the packet, and then releasing the associated buffer.

Packet transmission is often done at a lower priority than packet reception. This policy is superficially sound, because it minimizes the probability of packet loss when a burst of arriving packets exceeds the available buffer space. Reasonable operation of higher level protocols and applications, however, requires that transmit processing makes sufficient progress.

When the system is overloaded for long periods, use of a fixed lower priority for transmission leads to reduced throughput, or even complete cessation of packet transmission. Packets may be awaiting transmission, but the transmitting interface is idle. We call this *transmit starvation*.

Transmit starvation may occur if the transmitter interrupts at a lower priority than the receiver; or if they interrupt at the same priority, but the receiver's events are processed first by the driver; or if transmission completions are detected by polling, and the polling is done at a lower priority than receiver event processing.

This effect has also been described previously [17].

# 5. Avoiding livelock through better scheduling

In this section, we discuss several techniques to avoid receive livelocks. The techniques we discuss in this section include mechanisms to control the rate of incoming interrupts, polling-based mechanisms to ensure fair allocation of resources, and techniques to avoid unnecessary preemption.

## 5.1. Limiting the interrupt arrival rate

We can avoid or defer receive livelock by limiting the rate at which interrupts are imposed on the system. The system checks to see if interrupt processing is taking more than its share of resources, and if so, disables interrupts temporarily.

The system may infer impending livelock because it is discarding packets due to queue overflow, or because high-layer protocol processing or user-mode tasks are making no progress, or by measuring the fraction of CPU cycles used for packet processing. Once the system has invested enough work in an incoming packet to the point where it is about to be queued, it makes more sense to process that packet to completion than to drop it and rescue a subsequently-arriving packet from being dropped at the receiving interface, a cycle that could repeat *ad infinitum*.

When the system is about to drop a received packet because an internal queue is full, this strongly suggests that it should disable input interrupts. The host can then make progress on the packets already queued for higher-level processing, which has the side-effect of freeing buffers to use for subsequent received packets. Meanwhile, if the receiving interface has sufficient buffering of its own, additional incoming packets may accumulate there for a while.

We also need a trigger for re-enabling input interrupts, to prevent unnecessary packet loss. Interrupts may be re-enabled when internal buffer space becomes available, or upon expiration of a timer.

We may also want the system to guarantee some progress for user-level code. The system can observe that, over some interval, it has spent too much time processing packet input and output events, and temporarily disable interrupts to give higher protocol layers and user processes time to run. On a processor with a fine-grained clock register, the packet-input code can record the clock value on entry, subtract that from the clock value seen on exit, and keep a sum of the deltas. If this sum (or a running average) exceeds a specified fraction of the total elapsed time, the kernel disables input interrupts. (Digital's GIGAswitch™ system uses a similar mechanism [22].)

On a system without a fine-grained clock, one can crudely simulate this approach by sampling the CPU state on every clock interrupt (clock interrupts typically preempt device interrupt processing). If the system finds itself in the midst of processing interrupts for a series of such samples, it can disable interrupts for a few clock ticks.

## 5.2. Use of polling

Limiting the interrupt rate prevents system saturation but might not guarantee progress; the system must also fairly allocate packet-handling resources between input and output processing, and between multiple interfaces. We can provide fairness by carefully polling all sources of packet events, using a round-robin schedule.

In a pure polling system, the scheduler would invoke the device driver to ''listen'' for incoming packets and for transmit completion events. This would control the amount of device-level processing, and could also fairly allocate resources among event sources, thus avoiding livelock. Simply polling at fixed intervals, however, adds unacceptable latency to packet reception and transmission.

Polling designs and interrupt-driven designs differ in their placement of policy decisions. When the behavior of tasks cannot be predicted, we rely on the scheduler and the interrupt system to dynamically allocate CPU resources. When tasks can be expected to behave in a predictable manner, the tasks themselves are better able to make the scheduling decisions, and polling depends on voluntary cooperation among the tasks.

Since a purely interrupt-driven system leads to livelock, and a purely polling system adds unnecessary latency, we employ a hybrid design, in which the system polls only when triggered by an interrupt, and interrupts happen only while polling is suspended. During low loads, packet arrivals are unpredictable and we use interrupts to avoid latency. During high loads, we know that packets are arriving at or near the system's saturation rate, so we use polling to ensure progress and fairness, and only re-enable interrupts when no more work is pending.

## 5.3. Avoiding preemption

As we showed in section 4.2, receive livelock occurs because interrupt processing preempts all other packet processing. We can solve this problem by making higher-level packet processing non-preemptable. We observe that this can be done following one of two general approaches: do (almost) everything at high IPL, or do (almost) nothing at high IPL.

Following the first approach, we can modify the 4.2BSD design (see section 4.1) by eliminating the software interrupt, polling interfaces for events, and processing received packets to completion at device IPL. Because higher-level processing occurs at device IPL, it cannot be preempted by another packet arrival, and so we guarantee that livelock does not occur within the kernel's protocol stack. We still need to use a rate-control mechanism to ensure progress by user-level applications.

In a system following the second approach, the interrupt handler runs only long enough to set a ''service needed'' flag, and to schedule the polling thread if it is not already running. The

polling thread runs at zero IPL, checking the flags to decide which devices need service. Only when the polling thread is done does it re-enable the device interrupt. The polling thread can be interrupted at most once by each device, and so it progresses at full speed without interference.

Either approach eliminates the need to queue packets between the device driver and the higher-level protocol software, although if the protocol stack must block, the incoming packet must be queued at a later point. (For example, this would happen when the data is ready for delivery to a user process, or when an IP fragment is received and its companion fragments are not yet available.)

## 5.4. Summary of techniques

In summary, we avoid livelock by:

- Using interrupts only to initiate polling.

- Using round-robin polling to fairly allocate resources among event sources.

- Temporarily disabling input when feedback from a full queue, or a limit on CPU usage, indicates that other important tasks are pending.

- Dropping packets early, rather than late, to avoid wasted work. Once we decide to receive a packet, we try to process it to completion.

We maintain high performance by

- Re-enabling interrupts when no work is pending, to avoid polling overhead and to keep latency low.

- Letting the receiving interface buffer bursts, to avoid dropping packets.

- Eliminating the IP input queue, and associated overhead.

We observe, in passing, that inefficient code tends to exacerbate receive livelock, by lowering the MLFRR of the system and hence increasing the likelihood that livelock will occur. Aggressive optimization, ''fast-path'' designs, and removal of unnecessary steps all help to postpone arrival of livelock.

## 6. Livelock in BSD-based routers

In this section, we consider the specific example of an IP packet router built using Digital UNIX (formerly DEC OSF/1). We chose this application because routing performance is easily measured. Also, since firewalls typically use UNIX-based routers, they must be livelock-proof in order to prevent denial-of-service attacks.

Our goals were to (1) obtain the highest possible maximum throughput; (2) maintain high throughput even when overloaded; (3) allocate sufficient CPU cycles to user-mode tasks; (4) minimize latency; and (5) avoid degrading performance in other applications.

## 6.1. Measurement methodology

Our test configuration consisted of a router-under-test connecting two otherwise unloaded Ethernets. A source host generated IP/UDP packets at a variety of rates, and sent them via the router to a destination address. (The destination host did not exist; we fooled the router by inserting a phantom entry into its ARP table.) We measured router performance by counting the number of packets successfully forwarded in a given period, yielding an average forwarding rate.

The router-under-test was a DECstation™ 3000/300 Alpha-based system running Digital UNIX V3.2, with a SPECint92 rating of 66.2. We chose the slowest available Alpha host, to make the livelock problem more evident. The source host was a DECstation 3000/400, with a SPECint92 rating of 74.7. We slightly modified its kernel to allow more efficient generation of output packets, so that we could stress the router-under-test as much as possible.

In all the trials reported on here, the packet generator sent 10000 UDP packets carrying 4 bytes of data. This system does not generate a precisely paced stream of packets; the packet rates reported are averaged over several seconds, and the short-term rates varied somewhat from the mean. We calculated the delivered packet rate by using the ''netstat'' program (on the router machine) to sample the output interface count (''Opkts'') before and after each trial. We checked, using a network analyzer on the stub Ethernet, that this count exactly reports the number of packets transmitted on the output interface.

## 6.2. Measurements of an unmodified kernel

We started by measuring the performance of the unmodified operating system, as shown in figure 6-1. Each mark represents one trial. The filled circles show kernel-based forwarding performance, and the open squares show performance using the *screend* program [12], used in some firewalls to screen out unwanted packets. This user-mode program does one system call per packet; the packet-forwarding path includes both kernel and user-mode code. In this case, *screend* was configured to accept all packets.

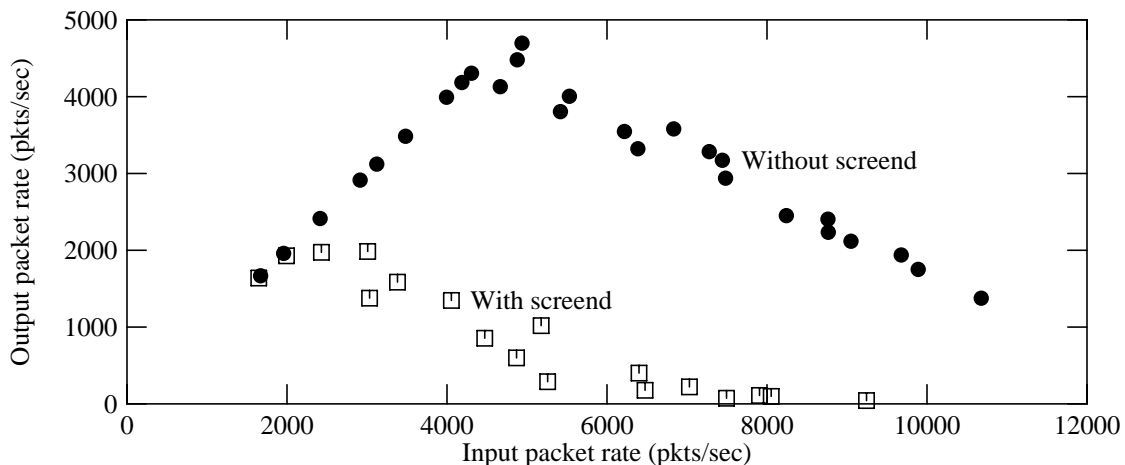

**Figure 6-1:** Forwarding performance of unmodified kernel

From these tests, it was clear that with *screend* running, the router suffered from poor overload behavior at rates above 2000 packets/sec., and complete livelock set in at about 6000

packets/sec. Even without *screend*, the router peaked at 4700 packets/sec., and would probably livelock somewhat below the maximum Ethernet packet rate of about 14,880 packets/second.

## 6.3. Why livelock occurs in the 4.2BSD model

4.2BSD follows the model described in section 4.1, and depicted in figure 6-2. The device driver runs at interrupt priority level (IPL) = SPLIMP, and the IP layer runs via a software interrupt at IPL = SPLNET, which is lower than SPLIMP. The queue between the driver and the IP code is named ''ipintrq,'' and each output interface is buffered by a queue of its own. All queues have length limits; excess packets are dropped. Device drivers in this system implement interrupt batching, so at high input rates very few interrupts are actually taken.

Digital UNIX follows a similar model, with the IP layer running as a separately scheduled thread at IPL = 0, instead of as a software interrupt handler.



**Figure 6-2:** IP forwarding path in 4.2BSD

It is now quite obvious why the system suffers from receive livelock. Once the input rate exceeds the rate at which the device driver can pull new packets out of the interface and add them to the IP input queue, the IP code never runs. Thus, it never removes packets from its queue (ipintrq), which fills up, and all subsequent received packets are dropped.

The system's CPU resources are saturated because it discards each packet after a lot of CPU time has been invested in it at elevated IPL. This is foolish; once a packet has made its way through the device driver, it represents an investment and should be processed to completion if at all possible. In a router, this means that the packet should be transmitted on the output interface. When the system is overloaded, it should discard packets as early as possible (i.e., in the receiving interface), so that discarded packets do not waste any resources.

## 6.4. Fixing the livelock problem

We solved the livelock problem by doing as much work as possible in a kernel thread, rather than in the interrupt handler, and by eliminating the IP input queue and its associated queue

manipulations and software interrupt (or thread dispatch)[1].  Once we decide to take a packet from the receiving interface, we try not to discard it later on, since this would represent wasted effort.

We also try to carefully ''schedule'' the work done in this thread.  It is probably not possible to use the system's real scheduler to control the handling of each packet, so we instead had this thread use a polling technique to efficiently simulate round-robin scheduling of packet processing.  The polling thread uses additional heuristics to help meet our performance goals.

In the new system, the interrupt handler for an interface driver does almost no work at all. Instead, it simple schedules the polling thread (if it has not already been scheduled), recording its need for packet processing, and then returns from the interrupt.  It does not set the device's interrupt-enable flag, so the system will not be distracted with additional interrupts until the polling thread has processed all of the pending packets.

At boot time, the modified interface drivers register themselves with the polling system, providing callback procedures for handling received and transmitted packets, and for enabling interrupts.  When the polling thread is scheduled, it checks all of the registered devices to see if they have requested processing, and invokes the appropriate callback procedures to do what the interrupt handler would have done in the unmodified kernel.

The received-packet callback procedures call the IP input processing routine directly, rather than placing received packets on a queue for later processing; this means that any packet accepted from the interface is processed as far as possible (e.g., to the output interface queue for forwarding, or to a queue for delivery to a process).  If the system falls behind, the interface's input buffer will soak up packets for a while, and any excess packets will be dropped by the interface before the system has wasted any resources on it.

The polling thread passes the callback procedures a quota on the number of packets they are allowed to handle.  Once a callback has used up its quota, it must return to the polling thread. This allows the thread to round-robin between multiple interfaces, and between input and output handling on any given interface, to prevent a single input stream from monopolizing the CPU.

Once all the packets pending at an interface have been handled, the polling thread also invokes the driver's interrupt-enable callback so that a subsequent packet event will cause an interrupt.

## 6.5. Results and analysis

Figures 6-3 summarizes the results of our changes, when *screend* is not used.  Several different kernel configurations are shown, using different mark symbols on the graph.  The modified kernel (shown with square marks) slightly improves the MLFRR, and avoids livelock at higher input rates.

---

[1]This is not such a radical idea; Van Jacobson had already used it as a way to improve end-system TCP performance [8].
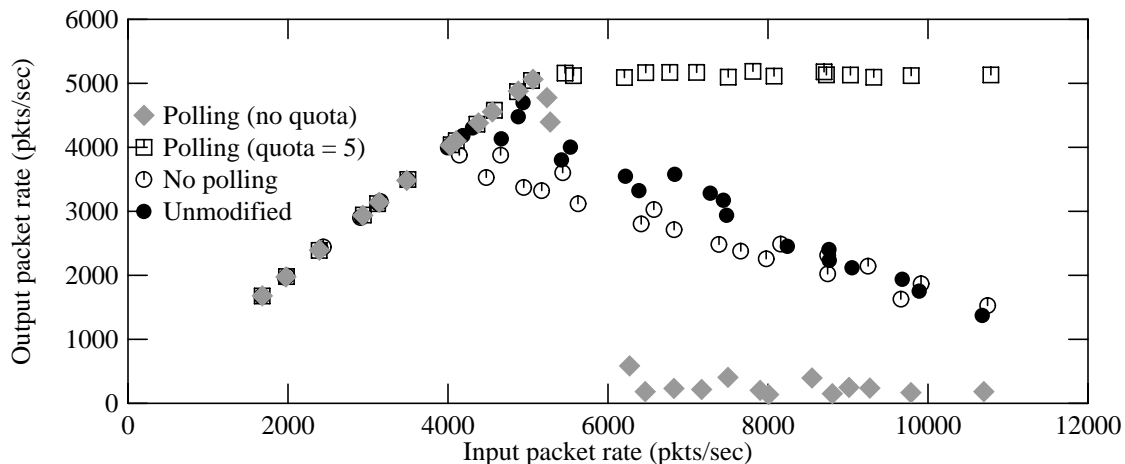
**Figure 6-3:** Forwarding performance of modified kernel, without using *screend*

The modified kernel can be configured to act as if it were an unmodified system (shown with open circles), although this seems to perform slightly worse than an actual unmodified system (filled circles). The reasons are not clear, but may involve slightly longer code paths, different compilers, or unfortunate changes in instruction cache conflicts.

## 6.6. Scheduling heuristics

Figure 6-3 shows that if the polling thread places no quota on the number of packets that a callback procedure can handle, when the input rate exceeds the MLFRR the total throughput drops almost to zero (shown with diamonds in the figure). This livelock occurs because although the packets are no longer discarded at the IP input queue, they are still piling up (and being discarded) at the queue for the output interface. This queue is unavoidable, since there is no guarantee that the output interface runs as fast as the input interface.

Why does the system fail to drain the output queue? If packets arrive too fast, the input-handling callback never finishes its job. This means that the polling thread never gets to call the output-handling callback for the transmitting interface, which prevents the release of transmitter buffer descriptors for use in further packet transmissions. This is similar to the transmit starvation condition identified in section 4.4.

The result is actually worse in the no-quota modified kernel, because in that system, packets are discarded for lack of space on the output queue, rather than on the IP input queue. The unmodified kernel does less work per discarded packet, and therefore occasionally discards them fast enough to catch up with a burst of input packets.

### 6.6.1. Feedback from full queues

How does the modified system perform when the *screend* program is used? Figure 6-4 compares the performance of the unmodified kernel (filled circles) and several modified kernels.

With the kernel modified as described so far (squares), the system performs about as badly as the unmodified kernel. The problem is that, because *screend* runs in user mode, the kernel must queue packets for delivery to *screend*. When the system is overloaded, this queue fills up and
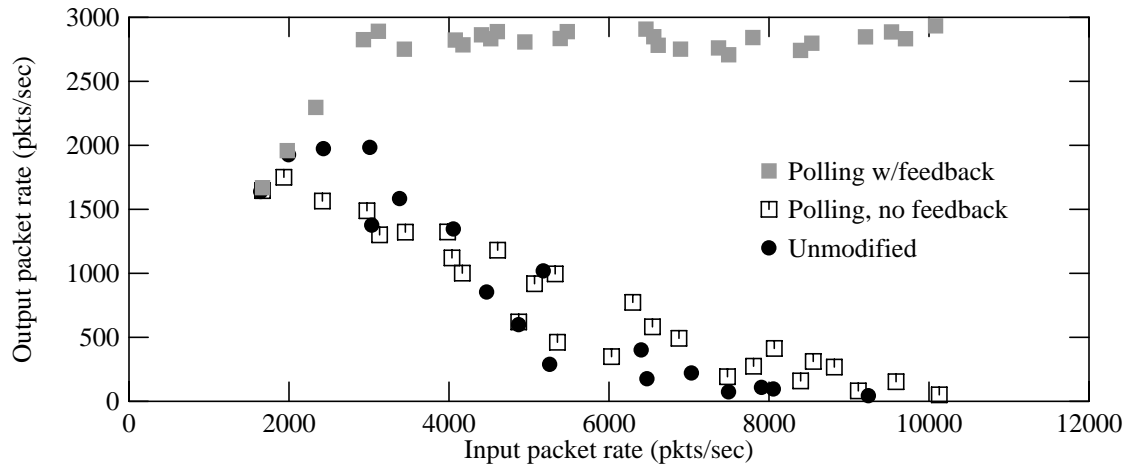
**Figure 6-4:** Forwarding performance of modified kernel, with *screend*

packets are dropped. *screend* never gets a chance to run to drain this queue, because the system devotes its cycles to handling input packets.

To resolve this problem, we detect when the screening queue becomes full and inhibit further input processing (and input interrupts) until more queue space is available. The result is shown with the gray square marks in figure 6-4: no livelock, and much improved peak throughput. Feedback from the queue state means that the system properly allocates CPU resources to move packets all the way through the system, instead of dropping them at an intermediate point.

In these experiments, the polling quota was 10 packets, the screening queue was limited to 32 packets, and we inhibited input processing when the queue was 75% full. Input processing is re-enabled when the screening queue becomes 25% full. We chose these high and low water marks arbitrarily, and some tuning might help. We also set a timeout (arbitrarily chosen as one clock tick, or about 1 msec) after which input is re-enabled, in case the *screend* program is hung, so that packets for other consumers are not dropped indefinitely.

The same queue-state feedback technique could be applied to other queues in the system, such as interface output queues, packet filter queues (for use in network monitoring) [14, 13], etc. The feedback policies for these queues would be more complex, since it might be difficult to determine if input processing load was actually preventing progress at these queues. Because the *screend* program is typically run as the only application on a system, however, a full screening queue is an unequivocal signal that too many packets are arriving.

### 6.6.2. Choice of packet-count quota

To avoid livelock in the non-*screend* configuration, we had to set a quota on the number of packets processed per callback, so we investigated how system throughput changes as the quota is varied. Figure 6-5 shows the results; smaller quotas work better. As the quota increases, livelock becomes more of a problem.

When *screend* is used, however, the queue-state feedback mechanism prevents livelock, and small quotas slightly reduce maximum throughput (by about 5%). We believe that by processing more packets per callback, the system amortizes the cost of polling more effectively, but increasing the quota could also increase worst-case per-packet latency. Once the quota is large enough

**Figure 6-5:** Effect of packet-count quota on performance, no *screend*

to fill the screening queue with a burst of packets, the feedback mechanism probably hides any potential for improvement.

Figure 6-6 shows the results when the *screend* process is in use.



**Figure 6-6:** Effect of packet-count quota on performance, with *screend*

In summary, tests both with and without *screend* suggest that a quota of between 5 and 10 packets yields stable and near-optimum behavior, for the hardware configuration tested. For other CPUs and network interfaces, the proper value may differ, so this parameter should be tunable.

# 7. Guaranteeing progress for user-level processes

The polling and queue-state feedback mechanisms described in section 6.4 can ensure that all necessary phases of packet processing make progress, even during input overload. They are indifferent to the needs of other activities, however, so user-level processes could still be starved for CPU cycles. This makes the system's user interface unresponsive and interferes with housekeeping tasks (such as routing table maintenance).

We verified this effect by running a compute-bound process on our modified router, and then flooding the router with minimum-sized packets to be forwarded. The router forwarded the packets at the full rate (i.e., as if no user-mode process were consuming resources), but the user process made no measurable progress.

Since the root problem is that the packet-input handling subsystem takes too much of the CPU, we should be able to ameliorate that by simply measuring the amount of CPU time spent handling received packets, and disabling input handling if this exceeds a threshold.

The Alpha architecture, on which we did these experiments, includes a high-resolution low-overhead counter register. This register counts every instruction cycle (in current implementations) and can be read in one instruction, without any data cache misses. Other modern RISC architectures support similar counters; Intel's Pentium is known to have one as an unsupported feature.

We measure the CPU usage over a period defined as several clock ticks (10 msec, in our current implementation, chosen arbitrarily to match the scheduler's quantum). Once each period, a timer function clears a running total of CPU cycles used in the packet-processing code.

Each time our modified kernel begins its polling loop, it reads the cycle counter, and reads it again at the end of the loop, to measure the number of cycles spent handling input and output packets during the loop. (The quota mechanism ensures that this interval is relatively short.) This number is then added to the running total, and if this total is above a threshold, input handling is immediately inhibited. At the end of the current period, a timer re-enables input handling. Execution of the system's idle thread also re-enables input interrupts and clears the running total.

By adjusting the threshold to be a fraction of the total number of cycles in a period, one can control fairly precisely the amount of CPU time spent processing packets. We have not yet implemented a programming interface for this control; for our tests, we simply patched a kernel global variable representing the percentage allocated to network processing, and the kernel automatically translates this to a number of cycles.
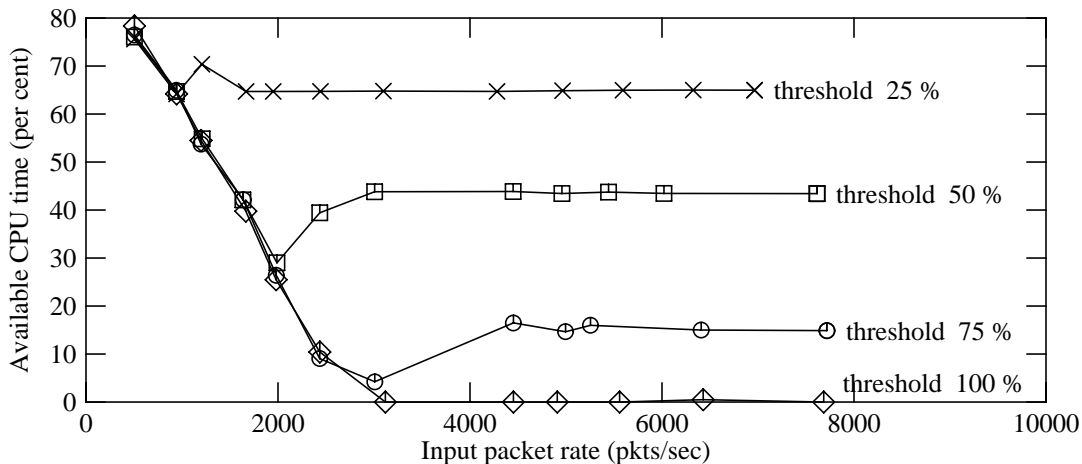


**Figure 7-1:** User-mode CPU time available using cycle-limit mechanism

Figure 7-1 shows how much CPU time is available to a compute-bound user process, for several settings of the cycle threshold and various input rates. The curves show fairly stable behavior as the input rate increases, but the user process does not get as much CPU time as the threshold setting would imply.

Part of the discrepancy comes from system overhead; even with no input load, the user process gets about 94% of the CPU cycles. Also, the cycle-limit mechanism inhibits packet input processing but not output processing. At higher input rates, before input is inhibited, the output queue fills enough to soak up additional CPU cycles.

Measurement error could cause some additional discrepancy. The cycle threshold is checked only after handling a burst of input packets (for these experiments, the callback quota was 5 packets). With the system forwarding about 5000 packets/second, handling such a burst takes about 1 msec, or about 10% of the threshold-checking period.

The initial dips in the curves for the 50% and 75% thresholds probably reflect the cost of handling the actual interrupts; these cycles are not counted against the threshold, and at input rates below saturation, each incoming packet may be handled fast enough that no interrupt batching occurs.

With a cycle-limit imposed on packet processing, the system is subjectively far more responsive, even during heavy input overload. This improvement, however, is mostly apparent for local users; any network-based interaction, such as Telnet, still suffers because many packets are being dropped.

## 7.1. Performance of end-system transport protocols

The changes we made to the kernel potentially affect the performance of end-system transport protocols, such as TCP and the UDP/RPC/XDR/NFS stack. Since we have not yet applied our modifications to a high-speed network interface driver, such as one for FDDI, we cannot yet measure this effect. (The test system can easily saturate an Ethernet, so measuring TCP throughput over Ethernet shows no effect.)

The technique of processing a received packet directly from the device driver to the TCP layer, without placing the packet on an IP-level queue, was used by Van Jacobson specifically to improve TCP performance [8]. It should reduce the cost of receiving a packet, by avoiding the queue operations and any associated locking; it also should improve the latency of kernel-to-kernel interactions (such as TCP acknowledgements and NFS RPCs).

The technique of polling the interfaces should not reduce end-system performance, because it is done primarily during input overload. (Some implementations use polling to avoid transmit interrupts altogether [10].) During overload, the unmodified system would not make any progress on applications or transport protocols; the use of polling, queue-state feedback, and CPU cycle limits should give the modified system a chance to make at least some progress.

Although TCP processing can be done in the same thread as the lower levels of the stack, NFS server implementations require separate threads of control, because a server may block waiting for disk I/O. Traditional NFS implementations, especially on servers, suffered badly from

livelock because during overload the NFS threads may never get a chance to run. With queue-state feedback from the NFS server's input queue, however, we should be able to avoid much of this problem. One could also use the CPU cycle-limit mechanism to reserve some resources for the NFS threads, although it might be difficult to find the ideal allocation.

NFS has typically used large UDP datagrams, and so causes frequent fragmentation of IP packets. Our kernel changes might reduce the tendency of the fragmentation reassembly mechanism to starve under transient overload conditions.

## 7.2. Measurements using traces of kernel execution

Although statistics showing the performance of a system under various loads enable one to compare the ultimate benefits of several approaches, these numbers do not provide a deep under-standing of the internal behavior of an operating system. We obtained traces of kernel execution to discover how the kernel is spending its time, and to measure the latency for several paths.

To obtain these traces, we used ATOM, an extremely flexible mechanism for instrumenting software [2, 3, 23]. ATOM takes a fully-linked binary program (even a Digital UNIX kernel) as input, and produces an instrumented binary as output. One also supplies to ATOM a module describing which points in the code to instrument, and a module containing analysis routines to execute at run-time.

Because ATOM allows the insertion of instrumentation at carefully chosen points in the ker-nel, it is possible to trace kernel paths without adding much overhead at all. On the DECstation 3000/300, a relatively slow Alpha system, tracing appeared to add about 1.5 usec per call or return. We did some trials in which almost all kernel procedures were traced, and others in which only a few were traced. The former trials provided insight into the precise code paths involved; the latter trials allowed us to obtain relatively accurate timing information.

Figure 7-2 shows most of the instrumentation code, except for some routines for determining which functions to instrument, and figure 7-3 shows all of the analysis code. The instrumen-tation module instructs ATOM to instrument calls to and returns from a selected set of procedures. The *ProcTrace()* function in the analysis module is then called at run-time, for each instrumented call or return, to record the current value of the cycle counter and a compact iden-tifier for the procedure involved.

An additional user-mode program (not shown) is used to extract the trace buffer from the kernel's memory, and to format it for later processing. This program is also used to reset the trace buffer at the beginning of a trial.

## 7.3. Traces of single-packet activity

We started by instrumenting almost all kernel procedures, except for a few low-level proce-dure that ATOM cannot currently trace and a small set of short but frequently invoked auxiliary procedures. Tables 7-1 and 7-2 briefly describe the procedures that appear in these traces.

We ran traces while using the system to forward a single minimum-length IP/UDP packet and extracted the relevant sequence of events. We could then plot these as timelines showing how

```
#include <stdio.h>
#include <string.h>
#include <cmplrs/atom.inst.h>
#include <assert.h>

/* Returns True if procedure should be instrumented */
int CanInstrument(Proc *p) {
    /* Code omitted. */
}

unsigned InstrumentAll() {       /* Called by ATOM */
  Obj *o; Proc *p;
  int index = 1;

  /* Prototype declarations for analysis routines */
  AddCallProto("Register(int, char*)");
  AddCallProto("ProcTrace(int, REGV)");

  o = GetFirstObj();
  if (BuildObj(o)) return 1;
  for (p = GetFirstObjProc(o); p != NULL; p = GetNextProc(p)) {
    if (CanInstrument(p)) {
      const char *name = ProcName(p);

      /* Register name of each instrumented routine,
                               with unique index */
      AddCallProgram(ProgramBefore, "Register", index, name);

      /* Instrument calls using index and
                              cycle-counter register */
      AddCallProc(p, ProcBefore, "ProcTrace", index,  REG_CC);
      /* Instrument returns using negated index */
      AddCallProc(p, ProcAfter, "ProcTrace", -index,  REG_CC);

      index += strlen(name)+1;
    }
  }
  WriteObj(o);
  return(0);
}
```

**Figure 7-2:** ATOM instrumentation code for kernel procedure tracing

procedure calls and returns nest, using the relative ''stack level'' to display the nesting. (Where the actual call stack includes uninstrumented procedures, the plotted stack level does not include calls through these procedures.)

Figure 7-4 shows a timeline for the modified kernel with polling disabled, which should approximate the behavior of an unmodified kernel. Figure 7-5 shows a timeline for the kernel with polling enabled. Each call is marked with the name of the procedure and the time at which the call was made, in microseconds since the start of the timeline. Returns are not individually marked, but one may deduce them from the decreases in stack level. Interrupts appear as if they were normal procedure calls.

In each case, we ran a rapid series of trials and selected one timeline in which no clock interrupts appear. To reduce the effects of cache misses, we never selected the first trial of a series. Even so, the times shown in these timelines should be treated as illustrative, but not necessarily typical. Also remember that instrumentation overhead adds several hundred microseconds to the total elapsed time (about 1.5 microseconds for each instrumented call or return).

```
#include "kptracets.h"  /* Defines SharedAtomData type */
SharedAtomData satom;

/*
 * Stores procedure name in buffer; caller uses buffer
 * offset as unique index value.
 */
void Register(int index, char *name) {
  char *pc = &satom.buffer[index];
  while(*pc++ = *name++)
    ;
}

/*
 * Called for each traced event.  Stores procedure index
 * and current cycle counter in parallel arrays.
 */
void ProcTrace(int index, int cycles) {
  if (satom.nextTrace < TRACE_SIZE) {
    satom.Procindex[satom.nextTrace] = index;
    satom.cycles[satom.nextTrace++] = cycles;
  }
}
```

**Figure 7-3:** ATOM analysis code for kernel procedure tracing

In figure 7-4, with polling disabled, we see the following interesting events (marked with dots on the timelines):

| | |
|---|---|
| 0 usec. | A packet has arrived, and lnintr() is called to begin handling the interrupt from the receiving LANCE Ethernet chip.  (Several microseconds have passed between interrupt assertion and the invocation of lnintr().) |
| 19 usec. | lnrint() is called to handle a received-packet interrupt. |
| 29 usec. | lnrint() calls lnread() to begin packet processing, which includes copying the packet to an mbuf structure. |
| 77 usec. | lnread() calls ether_input() to queue the received packet on the ipintr queue; ether_input() then calls netisr_input() to schedule a software interrupt. |
| 142 usec. | lnintr() finishes its execution at device IPL. |
| 191 usec. | After some thread-switching, ipinput() is invoked as a software interrupt. |
| 264 usec. | The IP-layer processing has determined that this packet should be forwarded, has chosen a next-hop destination, and now calls ip_output() to send the packet along. |
| 327 usec. | The LANCE driver has decided to send the packet, and calls lnput() to hand the buffer chain to the device. |
| 444 usec. | IP-layer processing is complete, and the software interrupt handler exits. |
| 522 usec. | The packet has been transmitted and the output interface has interrupted, causing a call to lnintr(). |
| 544 usec. | lntint() is called to handle the transmit interrupt. |
| 633 usec. | lntint() exits, completing all activity related to this packet. |

In figure 7-5, with polling enabled, we see a slightly different sequence of events:

| Procedure | Description |
|---|---|
| *Thread scheduling* | |
| thread_wakeup_prim | Used by thread scheduler to unblock a waiting thread. |
| thread_run | Switches between running threads. |
| assert_wait_mesg_head | Used by a thread to block on an event. |
| netisr_input | Notifies scheduler that the network software interrupt service routine should be running. |
| *Polling facility* | |
| lanpoll_isr* | Handler for software interrupt; polls devices with service requirements. |
| lanpoll_intsched* | Informs polling facility that an interrupt requires service. |
| *LANCE (Ethernet) driver* | |
| lnintr | Interrupt entry point for LANCE driver. |
| lnrint, lntint | Original (non-polling) receiver and transmitter interrupt service functions. |
| lnrintpoll*, lntintpoll* | New (for polling) receiver and transmitter interrupt service functions. |
| lnintena* | Called to re-enable LANCE interrupts. |
| lnread | Converts received packet buffer to mbuf chain. |
| lnoutput, lnstart | Initiates packet transmission. |
| lnput | Converts outgoing mbuf chain to packet buffer. |
| *Ethernet layer* | |
| ether_input | Parses MAC-level header of received packet. |
| ether_output | Adds MAC-level header to outgoing packet. |
| *IP layer* | |
| ipintr | Software interrupt handler for IP packet input. |
| ipinput | Parses IP header and dispatches received IP packet. |
| ip_output | Creates IP header for outgoing packet. |
| ip_forward | Forwards packets when host acts as a router. |
| *Clock interrupts* | |
| hardclock, clock_tick | Periodic (1024 Hz) clock interrupt handler |

*New routines added to support polling.

**Table 7-1:** Description of important procedures shown in timeline traces

0 usec.    A packet has arrived, and again lnintr() is called to begin handling the interrupt from the receiving LANCE chip.

21 usec.    lanpoll_intsched() is called to schedule a poll for this event.

53 usec.    lnintr() finishes its execution at device IPL. At this point, interrupts from this interface are still disabled, and the CPU is entirely under the control of the polling mechanism.

| Procedure | Description |
|---|---|
| thread_setrun,<br>thread_continue,<br>thread_block,<br>switch_context,<br>pmap_activate,<br>get_thread_high | Thread scheduling and memory management |
| malloc,<br>free | Memory allocation |
| m_leadingspace,<br>m_freem,<br>m_free,<br>m_copym | Mbuf manipulation |
| lninitdesc,<br>lnget | LANCE (Ethernet) driver |
| arpresolve_local | ARP layer |
| ip_forwardscreen,<br>in_canforward,<br>in_broadcast,<br>gw_forwardscreen | IP layer |
| bzero,<br>bcopy | Bulk memory operations |

**Table 7-2:** Description of boring procedures shown in timeline traces

97 usec.    After some thread-switching, lanpoll_isr() is called as a software interrupt handler, and begins its polling loop.

112 usec.   lnread() is called from lnrintpoll().

160 usec.   ether_input() determines that this is an IP packet, and does *not* place it on a queue.

166 usec.   ipinput() is called directly from ether_input().

235 usec.   The IP-layer processing calls ip_output() to send the packet along.

294 usec.   The LANCE driver calls lnput() to hand the buffer chain to the device.

407 usec.   IP-layer processing is complete, and control returns to the polling loop.

430 usec.   lanpoll_isr() calls lnintena() to re-enable interrupts from this device.

454 usec.   The packet has been transmitted and the output interface has interrupted, causing a call to lnintr(), which requests service for this event.

492 usec.   lanpoll_isr() is called without any thread-switching overhead, since this is still the current thread.

544 usec.   lntintpoll() is called to handle the transmit event.

586 usec.   lanpoll_isr() calls lnintena() to re-enable interrupts from this device.

597 usec.   lanpoll_isr() exits, completing all activity related to this packet.

From figures 7-4 and 7-5, one might conclude that with polling enabled, the kernel saves about 30 microseconds, mostly between the initial interrupt and the invocation of ipinput(). It is dangerous to base timing comparisons on a single pair of traces, and the instrumentation over-

**Figure 7-4:** Timeline forwarding a single packet, polling disabled

head confuses the situation somewhat, so we built another kernel just instrumenting the calls to lnintr() and ipinput(), and then ran a series of trials in order to obtain a statistically useful sample of the latency between these two points in the code. Each trial resulted in at least 10,000 packet receptions, almost all of which were short ICMP Echo packets.

The resulting distributions are shown in figure 7-6. The medians are marked with dots; the use of polling seems to reduce the median latency by about 13 usec. Polling also seems to reduce the latency variance somewhat. Note that the non-polling case includes the instrumentation overhead for one procedure return (from lnintr()) that is not included in the other case. The non-polling kernel also includes an extra "if" statement and an extra procedure call that were not present in the unmodified kernel, but these should not account for much time. In summary, we believe that the polling kernel, on this hardware, avoids about 10 microseconds of work per packet, probably because it does not move each packet onto and off off the ipintr queue.

Stack level

Stack level

lnintr 0

lanpoll_intsched 21
thread_wakeup_prim 26
thread_setrun 34

thread_run 61
pmap_activate 67
switch_context 73

thread_continue 84

lanpoll_isr 97
lnrintpoll 101

lnread 112

lnget 123
malloc 128

bcopy 145

ether_input 160
ipinput 166

ip_forwardscreen 184
new_forwardscreen 185
ip_forward 190
in_canforward 195

m_copym 205
malloc 209

bcopy 221

ip_output 235
in_broadcast 242

ether_output 254
arpresolve_local 259

m_leadingspace 272

lnoutput 284
lnstart 288
lnput 294
bcopy 299

bzero 314

bcopy 334

Time in usec

m_freem 368
m_free 374
free 379

lnintena 430

lnintr 454

lanpoll_intsched 473

lanpoll_isr 492
lntintpoll 497
bcopy 501

bcopy 511

m_freem 519
m_free 522
free 527

lninitdesc 540
bzero 546

bcopy 556

bcopy 569

lnstart 577

lnintena 586

Time in usec

**Figure 7-5:** Timeline forwarding a single packet, polling enabled

Cumulative fraction of events

87 usec
(Polling)

100 usec
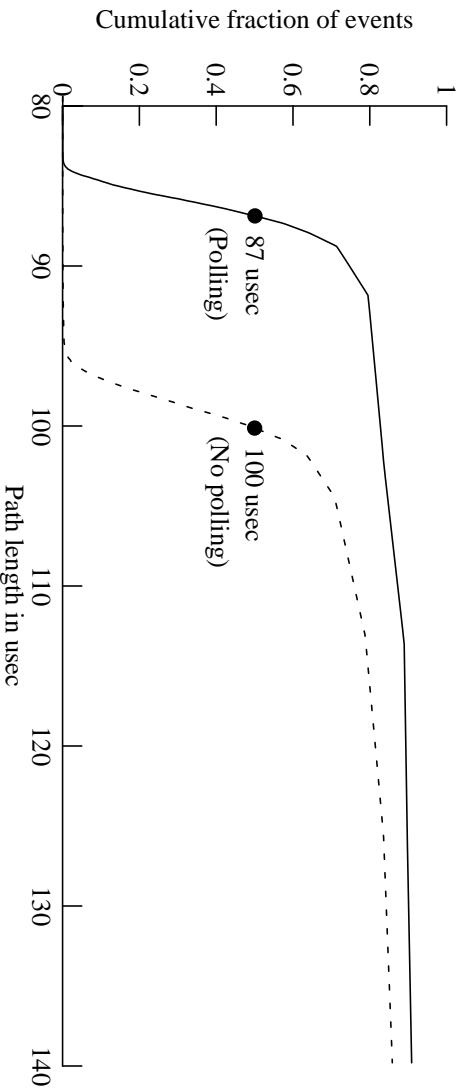(No polling)

Path length in usec

**Figure 7-6:** Distribution of latencies from lnintr() to ipinput()
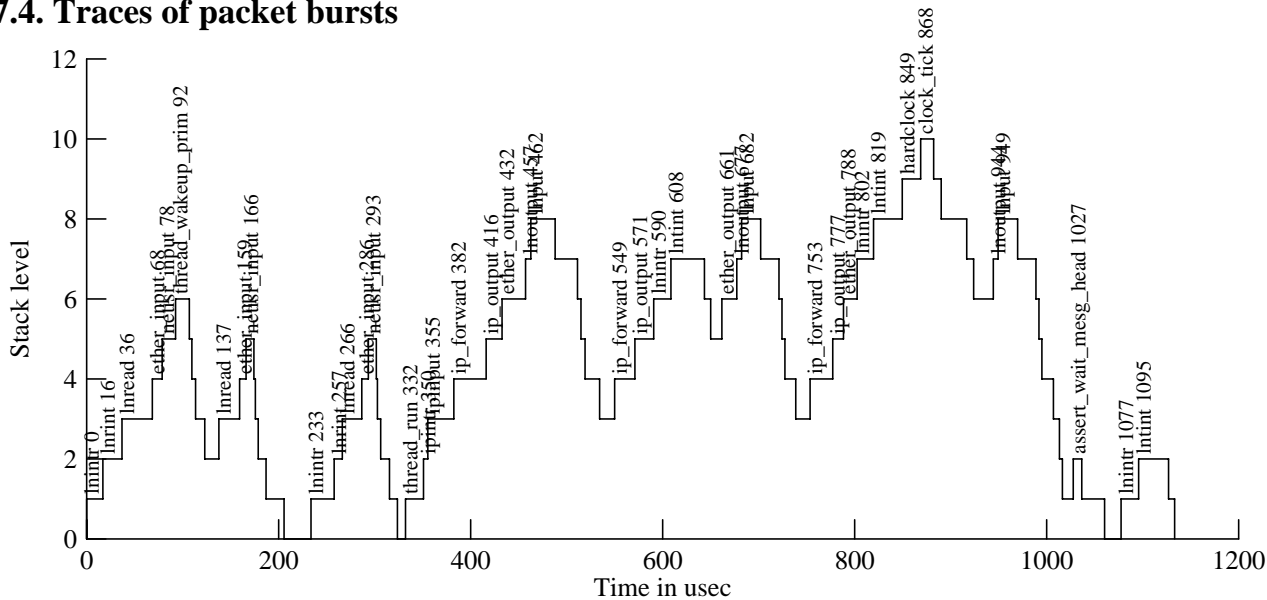
## 7.4. Traces of packet bursts



**Figure 7-7:** 3-packet burst latency, polling disabled

In section 4.3, we discussed how the unmodified kernel added extra latency to the processing of packets received in bursts. Figure 4-1 showed measurements of this effect on an ULTRIX kernel. With the ATOM tools, we can repeat this kind of measurement using our current test system. However, since ATOM cannot directly measure the assertion of the hardware interrupt signal, we do not include the kernel's initial interrupt latency. We believe this missing time amounts to less than 10 microseconds.

Figure 7-7 shows a traced timeline for the non-polling kernel handling a burst of 3 short packets to be forwarded. (We instrumented only the more interesting procedures for this trace, to reduce the clutter somewhat.) This trace starts with a call to lnintr(), which calls lnread() to queue the first packet, and then finds another waiting packet before calling lnread() again and dismissing the interrupt. Almost immediately, the third incoming packet causes another invocation of lnintr() and thus lnread().

The scheduler then invokes the software interrupt handler, which eventually calls ipinput(), which processes the first queued packet and eventually calls lnput() to place the forwarded packet in a buffer for transmission by the output interface. This happens 462 usec. after lnintr() is first called.

Figure 7-8 shows a similar timeline for the polling kernel. The first packet again results in a call to lnintr(), which schedules the polling thread and then immediately exists, leaving interrupts disabled for the receiving interface (and so the subsequent packet arrivals do not result in interrupts). The kernel then switches to the polling thread, which starts by calling back into the LANCE driver, and eventually gets to lnput() to place the first forwarded packet in a transmission buffer. This happens 366 microseconds after lnintr() is called, or almost 100 usec. sooner than with the non-polling kernel. Some part of this (perhaps 30 microseconds) comes from additional instrumentation overhead for the non-polling kernel.
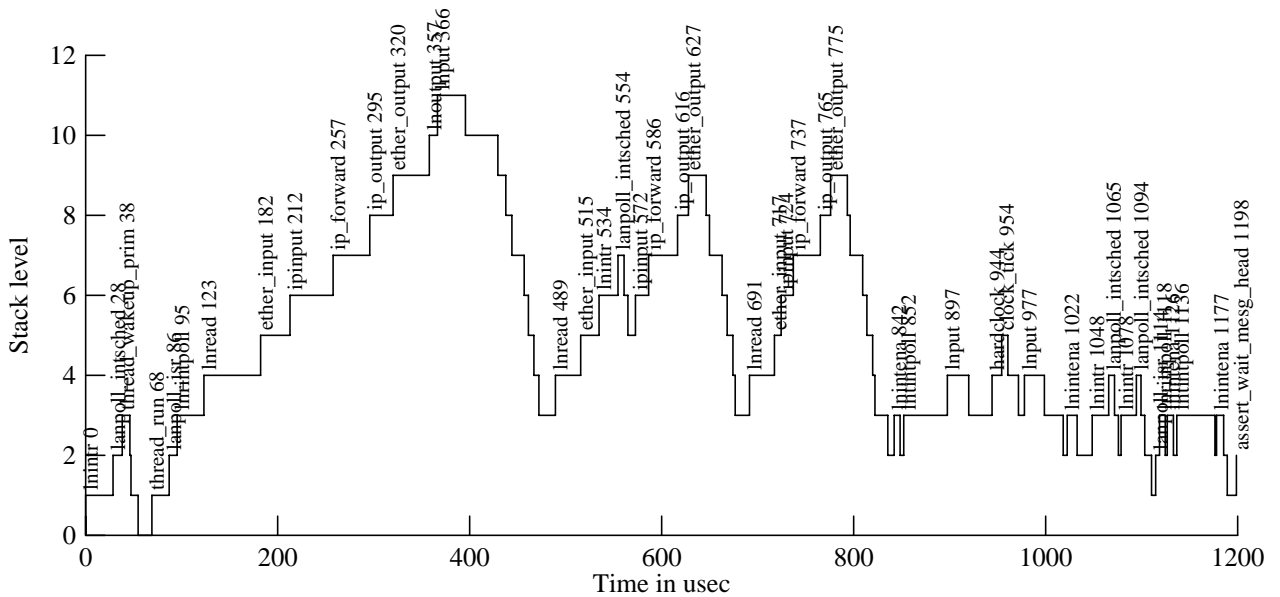
**Figure 7-8:** 3-packet burst latency, polling enabled

The difference may also be affected, in either direction, by some variation in the number of cache misses, although the trials were run repeatedly in order to warm up the caches. We also carefully selected traces that included as few clock interrupts as possible.

### 7.4.1. Implications for other applications

Although our changes improve end-to-end latency for the first packet in a burst of forwarded packets, packet forwarding is an unusual application because almost all of the work can be done without blocking. Most other applications, whether in-kernel (such as NFS service), or user-mode, require received packets to be queued for processing by another thread. Do our changes improve latency for these applications?

We note that as long as the polling thread has complete control of the CPU resources, nothing else can happen. In particular, no other thread can start processing the first packet of a burst. We can see two ways to avoid this problem:

- Multiprocessing: if the number of polling threads is smaller than the number of CPUs, at least some CPU resources will be available to finish processing early packets while the polling thread (or threads) continues to receive later packets.

- CPU-time limits: in section 7 we described how our polling mechanism can set a limit on the fraction of CPU time spent in the polling thread. We implemented this by disabling polling after the thread has used *m* milliseconds out of an *n* millisecond period. This has the side-effect of limiting first-packet latency (as seen by the next consumer after the polling thread) to approximately *m* milliseconds, unless the burst starts while polling is inhibited because of overload.

Neither of these completely solves the problem, but at least the polling mechanism provides these partial solutions; in a purely interrupt-driven kernel, user-mode application might have to wait for an unbounded interval to receive the first packet.

## 7.5. Possible improvements

The timeline in figure 7-8 suggests several ways in which the polling approach could be improved.

Although the polling kernel transmits the first packet of the burst sooner than the non-polling kernel, the latter has the advantage for the second packet, calling lnput() 682 microseconds after the start of the trace instead of 897 microseconds. The non-polling kernel is also marginally faster at sending the third packet (949 usec. instead of 977 usec.).

The cause for this discrepancy comes from an assumption in the code that until the transmitter interrupts, no new packets should be added to its output buffer. Until the interrupt is serviced, the interface is marked ''active'' and the upper-layer code leaves it alone. The non-polling kernel services the first transmitter interrupt (at 590 usec.) immediately, which allows it to restart the transmitter as soon as the second packet is ready for output. The polling kernel receives the interrupt sooner (at 534 usec.) but because the polling thread is still busy with the pending input packets, it fails to service the interrupt and so leaves the interface marked ''active.''

We believe that relatively straightforward changes would eliminate this extra latency. However, they would add some per-packet overhead: more frequent transmitter interrupts, and fewer chances to handle multiple output packets in one call to the driver layer. This overhead could reduce the MLFRR of the system, although it would not lead to livelock. It may be necessary to choose between optimizing MLFRR and optimizing latency.

The polling kernel also spends a little more time handling transmitter interrupts than the non-polling kernel, because after all three packets of the burst have been fully processed, the polling kernel still schedules the polling mechanism to see if anything else needs to be done. This does not add extra overhead during conditions of input overload, because then the polling mechanism would have useful work to do. However, at lower input rates it does rob cycles from other system tasks. We believe that there are several possible solutions to this problem, including different interrupt-generation schemes in the interface hardware, or some form of ''clocked interrupts'' [21, 24].

## 8. Avoiding livelock in a promiscuous network monitor

LAN monitoring applications typically require the host computer to place its network interface(s) into ''promiscuous mode,'' receiving all packets on the LAN, not just those addressed to the host itself. While a modern workstation can easily handle the full large-packet data rate of a high-speed LAN, if the LAN is flooded with small packets, even fast hosts might not keep up. For example, an FDDI LAN can carry up to 227,000 packets per second. At that rate, a host has about 4.4 usec. to process each packet.

## 8.1. Initial experiments

In our initial tests, we used a DECstation 3000/500 (SPECint92 = 74.3) running an unmodified DEC OSF/1 V3.0 kernel to monitor an FDDI LAN. The DECstation was running the *tcpdump* application [11], which was simply copying the first 40 bytes of each packet to `/dev/null`.

We used a packet generator to flood an FDDI LAN with 80-byte packets, and found that the monitoring system livelocked at between 21,000 and 24,000 packets/second At that point, interactive use was impossible. In figure 8-1, curve 1 shows the rate of packets delivered to the application; curve 2 shows the total rate received by the kernel, some of which were then dropped.
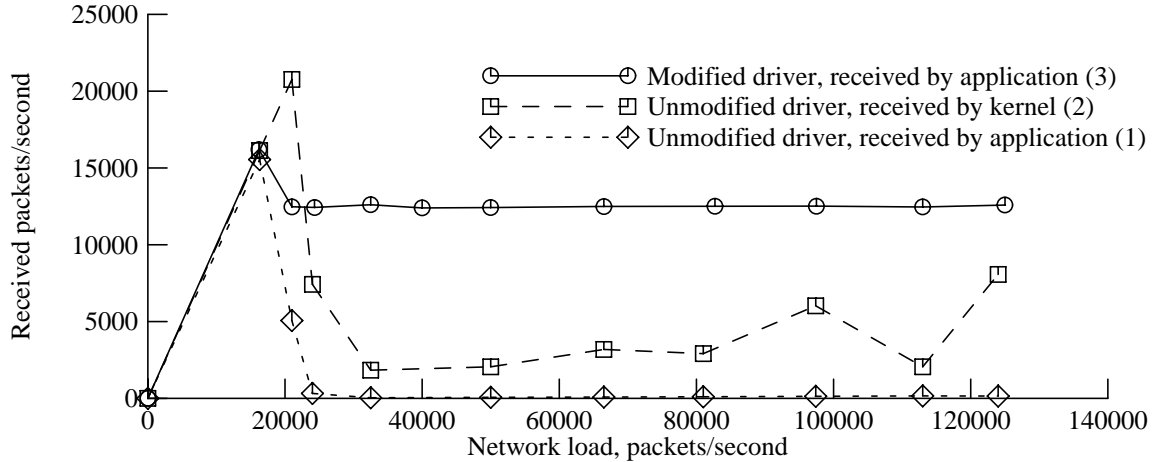


**Figure 8-1:** Network monitor performance, 80-byte FDDI packets

We then modified the FDDI receiver interrupt handler, which batches interrupts, to disable interrupts for 1 msec. once it had processed a batch of 250 packets. With this ''circuit-breaker'' change, response for all local (non-networked) interactive applications remained good regardless of the input load. The network monitor application continued to receive 100% of the input packets for loads up to at least 16,200 packets/sec., and then dropped to a nearly constant receive rate of about 12,400 packets/sec. as the network load increased to 125,000 packets/sec. (see curve 3 in figure 8-1).

Curve 3 shows that the reception rate during overload is somewhat less than the highest pre-overload rate. This is because our modified interrupt handler, although guaranteeing some progress to the application, does not directly avoid dropping the packets it receives. The circuit-breaker modification, while it avoids livelock, did not solve the problem as effectively as the polling kernel does.
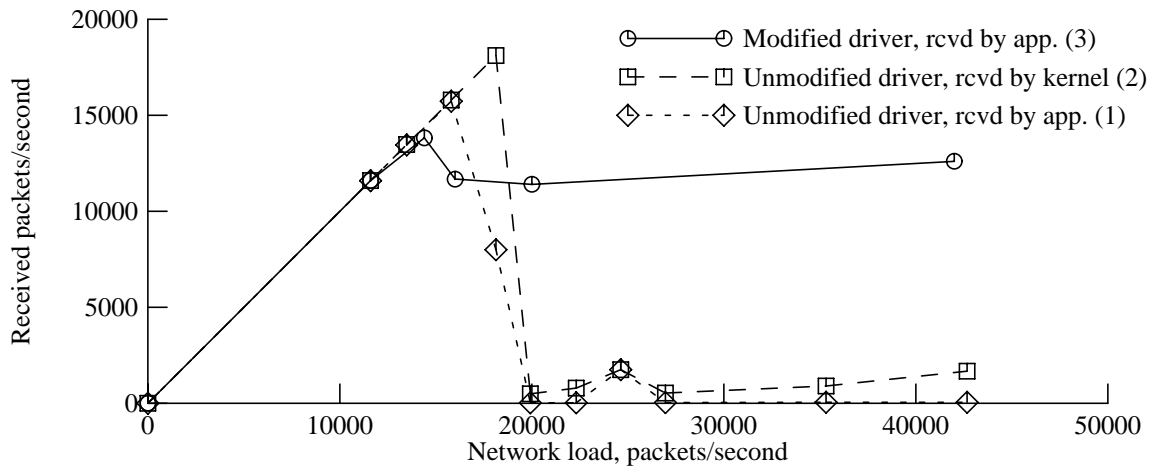


**Figure 8-2:** Network monitor performance, 272-byte FDDI packets

We also ran tests with larger (272-byte) FDDI packets, and found that the systems behaved about the same as they did with smaller packets: the unmodified system livelocked (curves 1 and 2, figure 8-2), and the system with the modified interface driver provided a nearly constant receive rate of about 12,000 packets/second. The similarity in MLFRR for the two packet sizes is unsurprising, because (except for some additional DMA) the kernel and application both do work independent of the packet size.

## 8.2. Network monitoring with the polling kernel

The *tcpdump* application obtains packets from the kernel using the packet filter pseudo-device driver [14]. Packet filtering (the selection of which received packets to hand to *tcpdump*) is done in the received-packet interrupt handler thread, and the resulting packets are put on a queue for delivery to the application. During overload conditions, the kernel discards packets because the application has no chance to drain this queue.

Note that *tcpdump* normally only looks at the packet headers, and so requests just the first 68 bytes of each packet. Since the kernel does not touch the remaining bytes of the packet, *tcpdump* throughput is nearly independent of packet size; its speed depends on per-packet, not per-byte, overheads.

We modified our polling kernel to implement queue-state feedback (see section 6.4) from the packet filter queue. For the measurements described here, the queue can contain at most 32 packets. Whenever the queue has room for fewer than 8 additional packets, we disable polling. We also set a 1 msec. timeout, after which polling is re-enabled. There is no direct mechanism to re-enable polling when more queue space has been made available, although perhaps there should be.

We then tested the network monitoring performance of the modified kernel. We set up a relatively slow system (a DECstation 3000/300) as the network monitor, and a faster system (a DECstation 3000/400) as a packet generator. The generator sent streams of packets to a third host on the same Ethernet, and *tcpdump* on the network monitor attempted to capture all of them, filtering on the UDP port number. For each trial, we sent between 10,000 and 30,000 packets at various rates, and measured the number of packets received by the monitor.

In all of these trials, we found that very few packets were dropped at the receiving interface. This means that almost all losses happened because the packet filter queue was full, not because the kernel failed to service the interface rapidly enough.

In our first set of trials, we configured *tcpdump* to simply copy the packet headers to the null device, /dev/null, instead of regular disk file. This should reduce the per-packet overhead and so increase the MLFRR. The results are shown in figure 8-3 for the polling kernel, with and without feedback.

The use of queue-state feedback clearly results in much better peak packet capture performance. It is tempting to infer that the use of feedback also improves overload behavior, but the rates measured in these trials do not quite reach the saturation point, and so provide no direct evidence about performance above that rate.
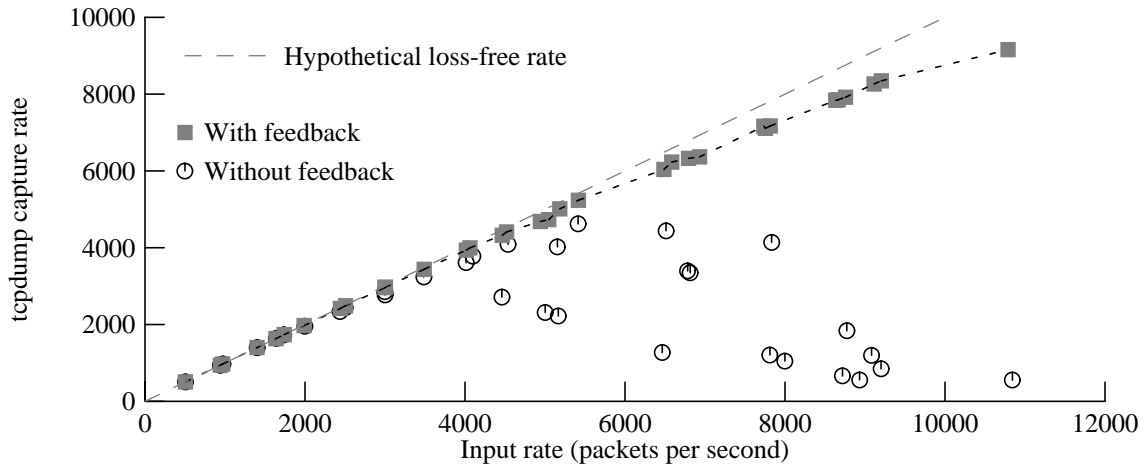
30

**Figure 8-3:** *tcpdump* capture rate, output to `/dev/null`

Note that at rates above the MLFRR, even with queue-state feedback the system does lose some packets. (The gray dashed line shows the performance of a hypothetical loss-free system.) We attribute this to the necessarily finite size of the packet filter queue: even though queue-state feedback inhibits input processing when the queue becomes full, the 1 msec. timeout happens before *tcpdump* has drained much of the queue, and so the kernel has no place to put the next batch of packets.

The results for the no-feedback kernel show a noisy relationship between input and output rates, above the MLFRR. This is because the packet generator is a relatively bursty source, and the mean burst size changes for different long-term generated rates. When the mean burst size is large, the network monitor processes more packets per interrupt, thus using fewer CPU cycles and leaving more for the *tcpdump* application.
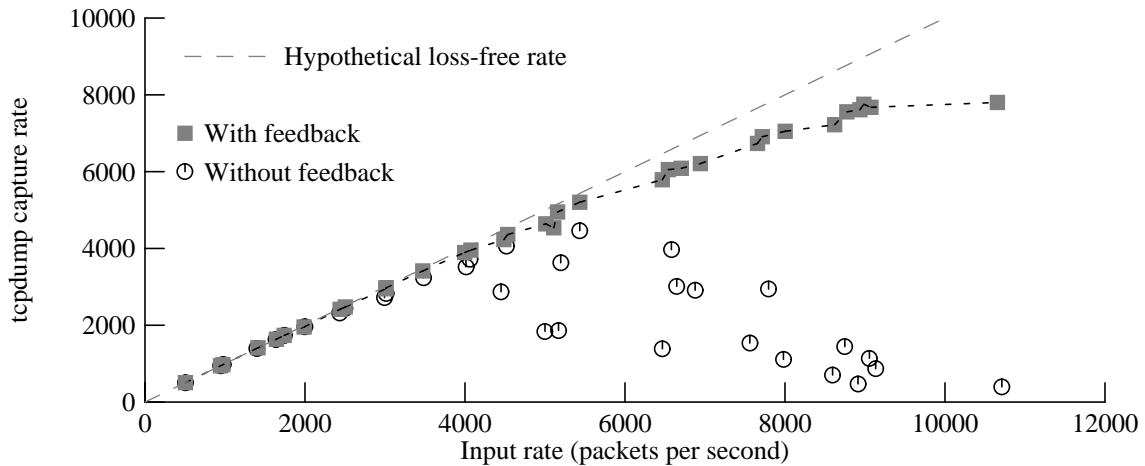


**Figure 8-4:** *tcpdump* capture rate, output to disk file

We then ran trials with *tcpdump* writing the received packet headers to a disk file (see figure 8-4). This added just enough per-packet overhead to allow us to saturate the system, even with queue-state feedback, at an input rate of about 9000 packets/sec (and a capture rate of about 7700 packets/sec).
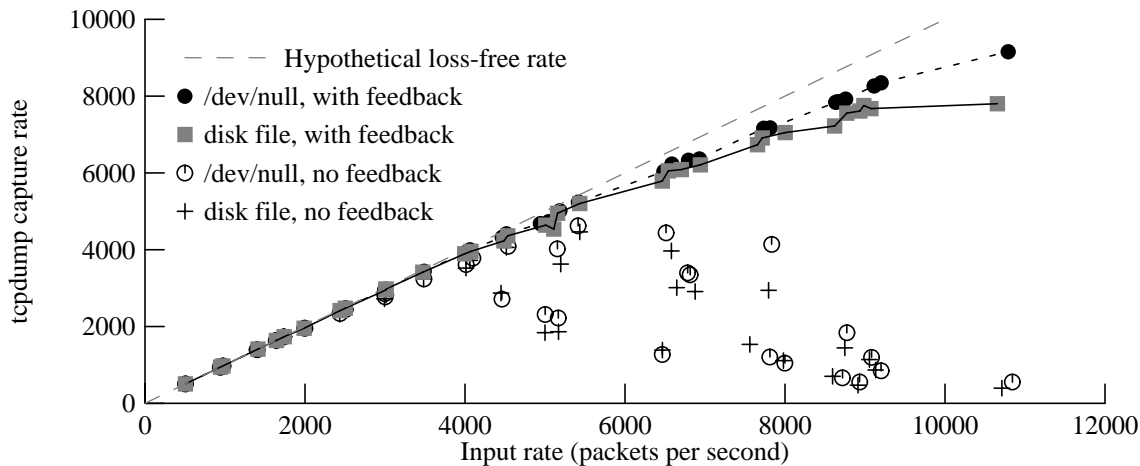
**Figure 8-5:** *tcpdump* capture rate, all trials shown

Figure 8-5 shows all four sets of trials together. From this, one can see that the extra overhead of writing packet headers to a disk file has at most a small effect on the capture rate, until the saturation point is reached.

# 9. Related work

Polling mechanisms have been used before in UNIX-based systems, both in network code and in other contexts. Whereas we have used polling to provide fairness and guaranteed progress, the previous applications of polling were intended to reduce the overhead associated with interrupt service. This does reduce the chances of system overload (for a given input rate), but does not prevent livelock.

Traw and Smith [21, 24] describe the use of ''clocked interrupts,'' periodic polling to learn of arriving packets without the overhead of per-packet interrupts. They point out that it is hard to choose the proper polling frequency: too high, and the system spends all its time polling; too low, and the receive latency soars. Their analysis [21] seems to ignore the use of interrupt batching to reduce the interrupt-service overhead; however, they do allude to the possibility of using a scheme in which an interrupt prompts polling for other events.

The 4.3BSD operating system [9] apparently used a periodic polling technique to process received characters from an eight-port terminal interface, if the recent input rate increased above a certain threshold. The intent seems to have been to avoid losing input characters (the device had little buffering available) but one could view this as a sort of livelock-avoidance strategy. Several router implementations use polling as their primary way to schedule packet processing.

When a congested router must drop a packet, its choice of which packet to drop can have significant effects. Our modifications do not affect *which* packets are dropped; we only change *when* they are dropped. The policy was and remains ''drop-tail''; other policies might provide better results [7].

Fall [4] discusses the problem of overload in non-flow-controlled systems such as routers. His approach improves system behavior by reducing per-packet and per-byte overheads, and thus increases the MLFRR, but does not directly improve overload behavior. His approach complements ours, but does not solve the livelock problem.

Some of our initial work on improved interface driver algorithms is described in [1].

## 9.1. Calaveras: scheduling in an embedded system's kernel

Some of the work presented in this paper was first done in the context of Calaveras, an advanced development project at Digital [18, 25]. The Calaveras kernel was designed to provide a foundation for building a high-bandwidth distributed file server capable of supporting traditional data as well as continuous media; in particular, it was used for a video-on-demand server. Preliminary experimental results with Calaveras showed that it did succeed in avoiding livelock.

The kernel distinguishes between three classes of tasks: isochronous, real-time, and general purpose. An isochronous task is a periodic time-driven task with performance requirements for bounded latency and jitter, and guaranteed throughput. A real-time task has performance requirements for low latency and high throughput. Both isochronous and real-time tasks have bounded execution times, so that the system can bound latency and jitter for other such tasks. In contrast, a general purpose task is characterized by lengthy or unbounded execution time, and may not have strict requirements for latency or throughput; however, the system does guarantee some progress for such tasks.

In the video-on-demand server, isochronous tasks were used primarily for servicing periodic activity such as video and audio transmission. Other activity, including non-isochronous I/O, was handled by real-time tasks.

The scheduler uses a Weighted Round Robin scheme for scheduling real-time and general-purpose tasks. Isochronous tasks are driven off timer interrupts.

The kernel assigns a scheduling flag and a weight to each real-time task. To schedule real-time tasks, the scheduler polls a set of flags, each of which indicates a pending task. For each flag that is set, the scheduler invokes the corresponding real-time task, with the assigned weight as a parameter. A scheduling flag may be set by an interrupt service routine as well as any real-time or general purpose task in the system.

The kernel introduced the notion of preemption windows. Isochronous and real-time tasks can only be preempted during their preemption windows. An isochronous task could be preempted by higher priority isochronous tasks; real-time tasks can be preempted by isochronous tasks. Preemption windows were placed at the completion of a unit of work, where little state needs to be saved in context-switching to the new task.

A real-time task is expected to process at most the number of work units equal to the weight passed to it as a parameter. Upon completion, the task saves its state, if necessary, and voluntarily returns to the scheduler.

After having completed one round of polling for real-time tasks, the scheduler switches to the list of general purpose tasks. This task runs until it blocks or the quantum timer expires. If the timer expires, the scheduler saves the tasks's context and goes back to servicing real-time tasks.

In this system, Ethernet and FDDI device drivers are implemented as real-time tasks. The only processing done at the interrupt level is to set a scheduling flag, to invoke the appropriate real-time task. General purpose tasks do not require any processing at the interrupt level.

The modified UNIX design we presented in section 6.4 corresponds quite closely to the design of the Calaveras scheduler. UNIX network driver processing corresponds to the Calaveras real-time tasks: a round-robin scheduler polls for a set of non-preemptable events, signalled but not scheduled by device interrupts. UNIX processes correspond to non-realtime tasks, which are preemptable and are scheduled at a lower priority.

The use of polling to schedule real-time tasks provides good performance during overload, and regulates the flow of traffic into the node. Just as flow control mechanisms such as a leaky bucket protect network resources from large bursts, polling protects the end-system resources by regulating the frequency at which work queues are polled, and by limiting the amount of work that may be performed during each poll.

Since each pending queue of tasks is serviced at a fixed interval, and task execution time is constant, or at least bounded, for all requests[2], we can easily provide delay bounds by limiting the length of each queue. Because interrupt service does not substantially preempt the processing of each packet, we can also maintain a reasonable bound on the latency of delivery of packets through the protocol stack.

## 10. Future work

Although the implementation described in this paper is straightforward and robust, and earlier versions have been deployed to customer installations, we see several areas that may require additional research.

### 10.1. Faster implementations

The experiments reported on in this paper were done on a relatively slow LAN (Ethernet, at 10 Mbits/sec.) and on the slowest available CPU that would run the Digital UNIX operating system. This allowed us to investigate the performance regime at and near overload, but the results cannot be extrapolated to predict the performance of state-of-the-art LANs and CPUs.

Use of faster CPUs in itself would be no problem; our kernel will run without modification on the fastest Alpha system. However, we would be unable to saturate such a system with our existing Ethernet-based test setup. To do high-speed experiments, we would have to obtain or construct a much faster packet generator. We would also have to apply our modifications to a driver for a faster LAN technology, such as FDDI or ATM. Unfortunately, such drivers seem to be far more complex than those for Ethernet interfaces.

Memory system performance has not usually improved quite as rapidly as the peak CPU instruction issue rate. This implies that future processors will be less sensitive to instruction counts, and more sensitive to memory locality. We expect that this will change the relative performance advantages of our kernel modifications, but we cannot predict which direction these changes will go.

_____

[2]An embedded system seldom needs to move data between different areas of memory, and therefore its processing delays may not depend on packet length. However, data checksums may cause length-dependent delays

## 10.2. Extension to multiprocessor kernels

Most computer vendors now sell some form of multiprocessor, to increase the performance of high-end systems beyond what is possible with a single CPU. Symmetric Multiprocessing (SMP) systems have been quite successful in many applications, and typically run a traditional operating system kernel that has been modified to support multiple kernel threads.

Although Digital UNIX is an SMP kernel, we have not yet extended our work to a true SMP environment. Our current kernel could run on SMP hardware, but it would do all of the interface driver and network processing on just one of the processors.

We believe that our polling approach allows better parallelization of interface and network processing, and so should improve performance on SMP systems. In order to demonstrate this, we would have to extend our modifications to include multiple polling threads, and we may have to modify the interface drivers to support concurrency in what were originally interrupt service routines (which are not run concurrently in the original system).

The use of multiple polling threads, while providing the opportunity for parallelization, also presents some challenges. How many threads should be active at once? If CPU cycle limits are enabled, should they apply to individually to each CPU, or to the entire system? And will the extra overhead of locking defeat the purpose of parallelization?

## 10.3. Selective packet dropping

Our approach to input overload is to drop packets as early as possible, to avoid spending resources on packets that would be dropped later on. This is a policy about when to drop packets, not about which packets to drop. In many cases, some packets may be much more important than others, and the system would be more effective if it preserved those when possible.

For example, when video streams are sent using Motion-JPEG compression, individual frames are independent of each other. A frame may be made up of several packets, all of which must be received for successful decompression. If one must discard several packets within a short interval, it is better to discard all of the packets of one frame rather than spreading the discarded packets among multiple frames, which then might all be impossible to decompress. Similarly, Romanow and Floyd have pointed out that if an ATM switch must drop multiple cells, it is better to drop all the cells of one packet rather than to spread the cell loss evenly over many packets [20].

Fall et al. have proposed *early discard load shedding* for Motion-JPEG streams, a technique in which frames are discarded as early in the processing pipeline as possible [5]. Their goal is to shed load before the system becomes overloaded, while we have focussed on responding to overload, but in either case it seems useful to be able to identify specific classes of packets to drop.

When received packets are dropped late, and therefore at upper levels in the network stack, it is fairly easy to determine which ones are useful and which ones are expendable. In our approach, because we have tried to drop packets as early as possible, perhaps before the software even sees them, it could be quite hard to drop all of the packets from one Motion-JPEG frame instead of spreading the losses at random.

Note that labelling techniques such as virtual circuits, ''flows,'' and packet priorities do not help here: from the sender's point of view, each Motion-JPEG frame has the same priority. One cannot *a priori* say that a given packet is more important than another; only when one is forced to drop a packet does it then become possible to define other packets as good targets for dropping.

## 10.4. Interactions with application-layer scheduling

We see several connections between scheduling the processing of received packets, and scheduling of application (user-mode) processes. At input rates below the MLFRR, it may be appropriate to modify the order in which packets are processed so as to reduce latency for packets destined to currently-running processes, or to provide batching across the kernel-user boundary (and so reduce context switching). At higher input rates, where some packets must be discarded, the packet-discard policy might favor packets for currently-scheduled processes; this could avoid thrashing. In either case, the process scheduler might prefer to give time to a process that has a large queue of pending packets.

Waldspurger has developed a proportional-share framework for expressing and implementing application scheduling policies [27, 28, 29]. In this framework, resource shares can be expressed in a kind of currency. He suggests [26] that if arriving packets, by the use of some simple labelling scheme, carry tokens in such a currency, it might be possible to integrate the packet-level scheduling decisions with process-level scheduling.

For example, suppose that packets were simply marked with the ID of the receiving process. As packets arrive for various processes, the kernel could track the number of unreceived packets per blocked process, and unblock the process with the largest batch of work to perform. If the system becomes overloaded, the network thread might start discarding all packets except for those destined to the current process.

Or suppose we want to drop excess packets according to some predetermined allocation of resources to processes. The network processing thread could keep track of how many packets have been received for each process during a recent interval, and discard packets in such a way as to maintain packet-consumption rates in proportion to the predetermined resource allocations.

We believe that such mechanisms could be implemented with minimal overhead, but we do not yet know how useful they would be.

## 11. Summary and conclusions

Systems that behave poorly under receive overload fail to provide consistent performance and good interactive behavior. Livelock is never the best response to overload. In this paper, we have shown how to understand system overload behavior and how to improve it, by carefully scheduling when packet processing is done.

We have shown, using measurements of a UNIX system, that traditional interrupt-driven systems perform badly under overload, resulting in receive livelock and starvation of transmits. Because such systems progressively reduce the priority of processing a packet as it goes further into

the system, when overloaded they exhibit excessive packet loss and wasted work. Such pathologies may be caused not only by long-term receive overload, but also by transient overload from short-term bursty arrivals.

We described a set of scheduling improvements that help solve the problem of poor overload behavior. These include:

- Limiting interrupt arrival rates, to shed overload

- Polling to provide fairness

- Processing received packets to completion

- Explicitly regulating CPU usage for packet processing

- Using feedback to inhibit input that would be discarded

Our experiments showed that these scheduling mechanisms provide good overload behavior and eliminate receive livelock. They should help both special-purpose and general-purpose systems.

## Acknowledgements

## References

[1]    Chran-Ham Chang, R. Flower, J. Forecast, H. Gray, W. R. Hawe, A. P Nadkarni, K. K. Ramakrishnan, U. N. Shikarpur, and K. M. Wilde. High-performance TCP/IP and UDP/IP Networking in DEC OSF/1 for Alpha AXP. *Digital Technical Journal* 5(1):44-61, Winter, 1993.

[2]    J. Bradley Chen and Alan Eustace. *Kernel Instrumentation Tools and Techniques*. Technical Report TR-26-95, Harvard University Center for Research in Computing Technology, November, 1995.

[3]    Alan Eustace and Amitabh Srivastava. ATOM: A Flexible Interface for Building High Performance Program Analysis Tools. In *Proc. 1995 USENIX Conference*, pages 303-313. New Orleans, January, 1995.

[4]    Kevin Fall. *A Peer-to-Peer I/O System in Support of I/O Intensive Workloads*. PhD thesis, University of California, San Diego, 1994.

[5]    Kevin Fall, Joseph Pasquale, and Steven McCanne. Workstation Video Playback Performance with Competitive Process Load. In *Proc. 5th International Workshop on Network and Operating Systems Support for Digital Audio and Video*, pages 179-182. Durham, NH, April, 1995.

[6]     Domenico Ferrari, Joseph Pasquale, and George C. Polyzos. *Network Issues for Sequoia 2000*. Sequoia 2000 Technical Report 91/6, University of California, Berkeley, December, 1991.

[7]     Sally Floyd and Van Jacobson.  Random Early Detection gateways for Congestion Avoidance. *Trans. Networking* 1(4):397-413, August, 1993.

[8]     Van Jacobson.  Efficient Protocol Implementation.  Notes from SIGCOMM '90 Tutorial on ''Protocols for High-Speed Networks''.  1990.

[9]     Samuel J. Leffler, Marshall Kirk McCusick, Michael J. Karels, and John S. Quarterman. *The Design and Implementation of the 4.3BSD UNIX Operating System.* Addison-Wesley, Reading, MA, 1989.

[10]    Rick Macklem.  Lessons Learned Tuning The 4.3BSD Reno Implementation of the NFS Protocol.  In *Proc. Winter 1991 USENIX Conference*, pages 53-64.  Dallas, TX, January, 1991.

[11]    Steven McCanne and Van Jacobson.  An Efficient, Extensible, and Portable Network Monitor. Work in progress.

[12]    Jeffrey C. Mogul.  Simple and Flexible Datagram Access Controls for Unix-based Gateways.  In *Proc. Summer 1989 USENIX Conference*, pages 203-221.  Baltimore, MD, June, 1989.

[13]    Jeffrey C. Mogul.  Efficient Use Of Workstations for Passive Monitoring of Local Area Networks.  In *Proc. SIGCOMM '90 Symposium on Communications Architectures and Protocols*, pages 253-263.  ACM SIGCOMM, Philadelphia, PA, September, 1990.

[14]    Jeffrey C. Mogul, Richard F. Rashid, Michael J. Accetta.  The Packet Filter:  An Efficient Mechanism for User-Level Network Code.  In *Proc. 11th Symposium on Operating Systems Principles*, pages 39-51.  Austin, Texas, November, 1987.

[15]    Radia Perlman.  Fault-Tolerant Broadcast of Routing Information. *Computer Networks* 7(6):395-405, December, 1983.

[16]    K. K. Ramakrishnan.  Scheduling Issues for Interfacing to High Speed Networks.  In *Proc. Globecom '92 IEEE Global Telecommunications Conf.*, pages 622-626.  Orlando, FL, December, 1992.

[17]    K. K. Ramakrishnan.  Performance Considerations in Designing Network Interfaces. *IEEE Journal on Selected Areas in Communications* 11(2):203-219, February, 1993.

[18]    K. K. Ramakrishnan, L. Vaitzblit, C. Gray, U. Vahalia, D. Ting, P. Tzelnic, S. Glaser, and W. Duso.  Operating System Support for a Video-on-Demand File Service. *ACM Multimedia Systems Journal* 3:53-65, March, 1995.

[19]    Marcus J. Ranum and Frederick M. Avolio.  A Toolkit and Methods for Internet Firewalls.  In *Proc. Summer 1994 USENIX Conference*, pages 37-44.  Boston, June, 1994.

[20]    Allyn Romanow and Sally Floyd.  Dynamics of TCP Traffic over ATM Networks. *IEEE J. Selected Areas in Communication* 13(4):633-641, May, 1995.

[21]    Jonathan M. Smith and C. Brendan S. Traw.  Giving Applications Access to Gb/s Networking. *IEEE Network* 7(4):44-52, July, 1993.

[22]    Robert J. Souza, P. G. Krishnakumar, Cüneyt M. Özveren, Robert J.Simcoe, Barry
        A. Spinney, Robert E. Thomas, and Robert J. Walsh.  GIGAswitch: A High-Performance
Packet Switching Platform.  *Digital Technical Journal* 6(1):9-22, Winter, 1994.

[23]    Amitabh Srivastava and Alan Eustace.  ATOM: A System for Building Customized
Program Analysis Tools.  In *Proc. SIGPLAN '94 Conf. on Programming Language Design and
Implementation*, pages 196-205.  Orlando, FL, June, 1994.

[24]    C. Brendan S. Traw and Jonathan M. Smith.  Hardware/Software Organization of a High-
Performance ATM Host Interface.  *IEEE Journal on Selected Areas in Communications*
11(2):240-253, February, 1993.

[25]    Uresh Vahalia, Cary G. Gray, and Dennis Ting.  Metadata Logging in an NFS Server.  In
*Proc. 1995 USENIX Conference*, pages 265-276.  New Orleans, LA, January, 1995.

[26]    Carl Waldspurger.  Private communication.

[27]    Carl A. Waldspurger and William E. Weihl.  Lottery Scheduling:  Flexible Proportional-
Share Resource Management.  In *Proc. First USENIX Symposium on Operating Systems Design
and Implementation (OSDI)*, pages 1-11.  Monterey, California, November, 1994.

[28]    Carl A. Waldspurger and William E. Weihl.  *Stride Scheduling:  Deterministic
Proportional-Share Resource Management*.  Technical Memorandum MIT/LCS/TM-528, Mas-
sachusetts Institute of Technology Laboratory for Computer Science, June, 1995.

[29]    Carl A. Waldspurger.  *Lottery and Stride Scheduling:  Flexible Proportional-Share
Resource Management*.  Technical Report (Ph.D. dissertation) MIT/LCS/TR-667, Massachusetts
Institute of Technology Laboratory for Computer Science, September, 1995.

# WRL Research Reports

''Titan System Manual.'' **Michael J. K. Nielsen.** WRL Research Report 86/1, September 1986.

''Global Register Allocation at Link Time.'' **David W. Wall.** WRL Research Report 86/3, October 1986.

''Optimal Finned Heat Sinks.'' **William R. Hamburgen.** WRL Research Report 86/4, October 1986.

''The Mahler Experience: Using an Intermediate Language as the Machine Description.'' **David W. Wall and Michael L. Powell.** WRL Research Report 87/1, August 1987.

''The Packet Filter: An Efficient Mechanism for User-level Network Code.'' **Jeffrey C. Mogul, Richard F. Rashid, Michael J. Accetta.** WRL Research Report 87/2, November 1987.

''Fragmentation Considered Harmful.'' **Christopher A. Kent, Jeffrey C. Mogul.** WRL Research Report 87/3, December 1987.

''Cache Coherence in Distributed Systems.'' **Christopher A. Kent.** WRL Research Report 87/4, December 1987.

''Register Windows vs. Register Allocation.'' **David W. Wall.** WRL Research Report 87/5, December 1987.

''Editing Graphical Objects Using Procedural Representations.'' **Paul J. Asente.** WRL Research Report 87/6, November 1987.

''The USENET Cookbook: an Experiment in Electronic Publication.'' **Brian K. Reid.** WRL Research Report 87/7, December 1987.

''MultiTitan: Four Architecture Papers.'' **Norman P. Jouppi, Jeremy Dion, David Boggs, Michael J. K. Nielsen.** WRL Research Report 87/8, April 1988.

''Fast Printed Circuit Board Routing.'' **Jeremy Dion.** WRL Research Report 88/1, March 1988.

''Compacting Garbage Collection with Ambiguous Roots.'' **Joel F. Bartlett.** WRL Research Report 88/2, February 1988.

''The Experimental Literature of The Internet: An Annotated Bibliography.'' **Jeffrey C. Mogul.** WRL Research Report 88/3, August 1988.

''Measured Capacity of an Ethernet: Myths and Reality.'' **David R. Boggs, Jeffrey C. Mogul, Christopher A. Kent.** WRL Research Report 88/4, September 1988.

''Visa Protocols for Controlling Inter-Organizational Datagram Flow: Extended Description.'' **Deborah Estrin, Jeffrey C. Mogul, Gene Tsudik, Kamaljit Anand.** WRL Research Report 88/5, December 1988.

''SCHEME->C A Portable Scheme-to-C Compiler.'' **Joel F. Bartlett.** WRL Research Report 89/1, January 1989.

''Optimal Group Distribution in Carry-Skip Adders.'' **Silvio Turrini.** WRL Research Report 89/2, February 1989.

''Precise Robotic Paste Dot Dispensing.'' **William R. Hamburgen.** WRL Research Report 89/3, February 1989.

''Simple and Flexible Datagram Access Controls for Unix-based Gateways.'' **Jeffrey C. Mogul.** WRL Research Report 89/4, March 1989.

''Spritely NFS: Implementation and Performance of Cache-Consistency Protocols.'' **V. Srinivasan and Jeffrey C. Mogul.** WRL Research Report 89/5, May 1989.

''Available Instruction-Level Parallelism for Superscalar and Superpipelined Machines.'' **Norman P. Jouppi and David W. Wall.** WRL Research Report 89/7, July 1989.

''A Unified Vector/Scalar Floating-Point Architecture.'' **Norman P. Jouppi, Jonathan Bertoni, and David W. Wall.** WRL Research Report 89/8, July 1989.

''Architectural and Organizational Tradeoffs in the Design of the MultiTitan CPU.'' **Norman P. Jouppi.** WRL Research Report 89/9, July 1989.

''Integration and Packaging Plateaus of Processor Performance.'' **Norman P. Jouppi.** WRL Research Report 89/10, July 1989.

''A 20-MIPS Sustained 32-bit CMOS Microprocessor with High Ratio of Sustained to Peak Performance.'' **Norman P. Jouppi and Jeffrey Y. F. Tang.** WRL Research Report 89/11, July 1989.

''The Distribution of Instruction-Level and Machine Parallelism and Its Effect on Performance.'' **Norman P. Jouppi.** WRL Research Report 89/13, July 1989.

''Long Address Traces from RISC Machines: Generation and Analysis.'' **Anita Borg, R.E.Kessler, Georgia Lazana, and David W. Wall.** WRL Research Report 89/14, September 1989.

''Link-Time Code Modification.'' **David W. Wall.** WRL Research Report 89/17, September 1989.

''Noise Issues in the ECL Circuit Family.'' **Jeffrey Y.F. Tang and J. Leon Yang.** WRL Research Report 90/1, January 1990.

''Efficient Generation of Test Patterns Using Boolean Satisfiablilty.'' **Tracy Larrabee.** WRL Research Report 90/2, February 1990.

''Two Papers on Test Pattern Generation.'' **Tracy Larrabee.** WRL Research Report 90/3, March 1990.

''Virtual Memory vs. The File System.'' **Michael N. Nelson.** WRL Research Report 90/4, March 1990.

''Efficient Use of Workstations for Passive Monitoring of Local Area Networks.'' **Jeffrey C. Mogul.** WRL Research Report 90/5, July 1990.

''A One-Dimensional Thermal Model for the VAX 9000 Multi Chip Units.'' **John S. Fitch.** WRL Research Report 90/6, July 1990.

''1990 DECWRL/Livermore Magic Release.'' **Robert N. Mayo, Michael H. Arnold, Walter S. Scott, Don Stark, Gordon T. Hamachi.** WRL Research Report 90/7, September 1990.

''Pool Boiling Enhancement Techniques for Water at Low Pressure.'' **Wade R. McGillis, John S. Fitch, William R. Hamburgen, Van P. Carey.** WRL Research Report 90/9, December 1990.

''Writing Fast X Servers for Dumb Color Frame Buffers.'' **Joel McCormack.** WRL Research Report 91/1, February 1991.

''A Simulation Based Study of TLB Performance.'' **J. Bradley Chen, Anita Borg, Norman P. Jouppi.** WRL Research Report 91/2, November 1991.

''Analysis of Power Supply Networks in VLSI Circuits.'' **Don Stark.** WRL Research Report 91/3, April 1991.

''TurboChannel T1 Adapter.'' **David Boggs.** WRL Research Report 91/4, April 1991.

''Procedure Merging with Instruction Caches.'' **Scott McFarling.** WRL Research Report 91/5, March 1991.

''Don't Fidget with Widgets, Draw!.'' **Joel Bartlett.** WRL Research Report 91/6, May 1991.

''Pool Boiling on Small Heat Dissipating Elements in Water at Subatmospheric Pressure.'' **Wade R. McGillis, John S. Fitch, William R. Hamburgen, Van P. Carey.** WRL Research Report 91/7, June 1991.

''Incremental, Generational Mostly-Copying Garbage Collection in Uncooperative Environments.'' **G. May Yip.** WRL Research Report 91/8, June 1991.

''Interleaved Fin Thermal Connectors for Multichip Modules.'' **William R. Hamburgen.** WRL Research Report 91/9, August 1991.

''Experience with a Software-defined Machine Architecture.'' **David W. Wall.** WRL Research Report 91/10, August 1991.

''Network Locality at the Scale of Processes.'' **Jeffrey C. Mogul.** WRL Research Report 91/11, November 1991.

''Cache Write Policies and Performance.'' **Norman P. Jouppi.** WRL Research Report 91/12, December 1991.

''Packaging a 150 W Bipolar ECL Microprocessor.'' **William R. Hamburgen, John S. Fitch.** WRL Research Report 92/1, March 1992.

''Observing TCP Dynamics in Real Networks.'' **Jeffrey C. Mogul.** WRL Research Report 92/2, April 1992.

''Systems for Late Code Modification.'' **David W. Wall.** WRL Research Report 92/3, May 1992.

''Piecewise Linear Models for Switch-Level Simulation.'' **Russell Kao.** WRL Research Report 92/5, September 1992.

''A Practical System for Intermodule Code Optimization at Link-Time.'' **Amitabh Srivastava and David W. Wall.** WRL Research Report 92/6, December 1992.

''A Smart Frame Buffer.'' **Joel McCormack & Bob McNamara.** WRL Research Report 93/1, January 1993.

''Recovery in Spritely NFS.'' **Jeffrey C. Mogul.** WRL Research Report 93/2, June 1993.

''Tradeoffs in Two-Level On-Chip Caching.'' **Norman P. Jouppi & Steven J.E. Wilton.** WRL Research Report 93/3, October 1993.

''Unreachable Procedures in Object-oriented Programing.'' **Amitabh Srivastava.** WRL Research Report 93/4, August 1993.

''An Enhanced Access and Cycle Time Model for On-Chip Caches.'' **Steven J.E. Wilton and Norman P. Jouppi.** WRL Research Report 93/5, July 1994.

''Limits of Instruction-Level Parallelism.'' **David W. Wall.** WRL Research Report 93/6, November 1993.

''Fluoroelastomer Pressure Pad Design for Microelectronic Applications.'' **Alberto Makino, William R. Hamburgen, John S. Fitch.** WRL Research Report 93/7, November 1993.

''A 300MHz 115W 32b Bipolar ECL Microprocessor.'' **Norman P. Jouppi, Patrick Boyle, Jeremy Dion, Mary Jo Doherty, Alan Eustace, Ramsey Haddad, Robert Mayo, Suresh Menon, Louis Monier, Don Stark, Silvio Turrini, Leon Yang, John Fitch, William Hamburgen, Russell Kao, and Richard Swan.** WRL Research Report 93/8, December 1993.

''Link-Time Optimization of Address Calculation on a 64-bit Architecture.'' **Amitabh Srivastava, David W. Wall.** WRL Research Report 94/1, February 1994.

''ATOM: A System for Building Customized Program Analysis Tools.'' **Amitabh Srivastava, Alan Eustace.** WRL Research Report 94/2, March 1994.

''Complexity/Performance Tradeoffs with Non-Blocking Loads.'' **Keith I. Farkas, Norman P. Jouppi.** WRL Research Report 94/3, March 1994.

''A Better Update Policy.'' **Jeffrey C. Mogul.** WRL Research Report 94/4, April 1994.

''Boolean Matching for Full-Custom ECL Gates.'' **Robert N. Mayo, Herve Touati.** WRL Research Report 94/5, April 1994.

''Software Methods for System Address Tracing: Implementation and Validation.'' **J. Bradley Chen, David W. Wall, and Anita Borg.** WRL Research Report 94/6, September 1994.

''Performance Implications of Multiple Pointer Sizes.'' **Jeffrey C. Mogul, Joel F. Bartlett, Robert N. Mayo, and Amitabh Srivastava.** WRL Research Report 94/7, December 1994.

''How Useful Are Non-blocking Loads, Stream Buffers, and Speculative Execution in Multiple Issue Processors?.'' **Keith I. Farkas, Norman P. Jouppi, and Paul Chow.** WRL Research Report 94/8, December 1994.

''Drip: A Schematic Drawing Interpreter.'' **Ramsey W. Haddad.** WRL Research Report 95/1, March 1995.

''Recursive Layout Generation.'' **Louis M. Monier, Jeremy Dion.** WRL Research Report 95/2, March 1995.

''Contour: A Tile-based Gridless Router.'' **Jeremy Dion, Louis M. Monier.** WRL Research Report 95/3, March 1995.

''The Case for Persistent-Connection HTTP.'' **Jeffrey C. Mogul.** WRL Research Report 95/4, May 1995.

''Network Behavior of a Busy Web Server and its Clients.'' **Jeffrey C. Mogul.** WRL Research Report 95/5, October 1995.

''The Predictability of Branches in Libraries.'' **Brad Calder, Dirk Grunwald, and Amitabh Srivastava.** WRL Research Report 95/6, October 1995.

''Shared Memory Consistency Models: A Tutorial.'' **Sarita V. Adve, Kourosh Gharachorloo.** WRL Research Report 95/7, September 1995.

''Eliminating Receive Livelock in an Interrupt-driven Kernel.'' **Jeffrey C. Mogul and K. K. Ramakrishnan.** WRL Research Report 95/8, December 1995.

# WRL Technical Notes

''TCP/IP PrintServer: Print Server Protocol.'' **Brian K. Reid and Christopher A. Kent.** WRL Technical Note TN-4, September 1988.

''TCP/IP PrintServer: Server Architecture and Implementation.'' **Christopher A. Kent.** WRL Technical Note TN-7, November 1988.

''Smart Code, Stupid Memory: A Fast X Server for a Dumb Color Frame Buffer.'' **Joel McCormack.** WRL Technical Note TN-9, September 1989.

''Why Aren't Operating Systems Getting Faster As Fast As Hardware?.'' **John Ousterhout.** WRL Technical Note TN-11, October 1989.

''Mostly-Copying Garbage Collection Picks Up Generations and C++.'' **Joel F. Bartlett.** WRL Technical Note TN-12, October 1989.

''Characterization of Organic Illumination Systems.'' **Bill Hamburgen, Jeff Mogul, Brian Reid, Alan Eustace, Richard Swan, Mary Jo Doherty, and Joel Bartlett.** WRL Technical Note TN-13, April 1989.

''Improving Direct-Mapped Cache Performance by the Addition of a Small Fully-Associative Cache and Prefetch Buffers.'' **Norman P. Jouppi.** WRL Technical Note TN-14, March 1990.

''Limits of Instruction-Level Parallelism.'' **David W. Wall.** WRL Technical Note TN-15, December 1990.

''The Effect of Context Switches on Cache Performance.'' **Jeffrey C. Mogul and Anita Borg.** WRL Technical Note TN-16, December 1990.

''MTOOL: A Method For Detecting Memory Bottlenecks.'' **Aaron Goldberg and John Hennessy.** WRL Technical Note TN-17, December 1990.

''Predicting Program Behavior Using Real or Estimated Profiles.'' **David W. Wall.** WRL Technical Note TN-18, December 1990.

''Cache Replacement with Dynamic Exclusion.'' **Scott McFarling.** WRL Technical Note TN-22, November 1991.

''Boiling Binary Mixtures at Subatmospheric Pressures.'' **Wade R. McGillis, John S. Fitch, William R. Hamburgen, Van P. Carey.** WRL Technical Note TN-23, January 1992.

''A Comparison of Acoustic and Infrared Inspection Techniques for Die Attach.'' **John S. Fitch.** WRL Technical Note TN-24, January 1992.

''TurboChannel Versatec Adapter.'' **David Boggs.** WRL Technical Note TN-26, January 1992.

''A Recovery Protocol For Spritely NFS.'' **Jeffrey C. Mogul.** WRL Technical Note TN-27, April 1992.

''Electrical Evaluation Of The BIPS-0 Package.'' **Patrick D. Boyle.** WRL Technical Note TN-29, July 1992.

''Transparent Controls for Interactive Graphics.'' **Joel F. Bartlett.** WRL Technical Note TN-30, July 1992.

''Design Tools for BIPS-0.'' **Jeremy Dion & Louis Monier.** WRL Technical Note TN-32, December 1992.

''Link-Time Optimization of Address Calculation on a 64-Bit Architecture.'' **Amitabh Srivastava and David W. Wall.** WRL Technical Note TN-35, June 1993.

''Combining Branch Predictors.'' **Scott McFarling.** WRL Technical Note TN-36, June 1993.

''Boolean Matching for Full-Custom ECL Gates.'' **Robert N. Mayo and Herve Touati.** WRL Technical Note TN-37, June 1993.

''Piecewise Linear Models for Rsim.'' **Russell Kao, Mark Horowitz.** WRL Technical Note TN-40, December 1993.

''Speculative Execution and Instruction-Level Parallelism.'' **David W. Wall.** WRL Technical Note TN-42, March 1994.

''Ramonamap - An Example of Graphical Groupware.'' **Joel F. Bartlett.** WRL Technical Note TN-43, December 1994.

''ATOM: A Flexible Interface for Building High Performance Program Analysis Tools.'' **Alan Eustace and Amitabh Srivastava.** WRL Technical Note TN-44, July 1994.

''Circuit and Process Directions for Low-Voltage Swing Submicron BiCMOS.'' **Norman P. Jouppi, Suresh Menon, and Stefanos Sidiropoulos.** WRL Technical Note TN-45, March 1994.

''Experience with a Wireless World Wide Web Client.'' **Joel F. Bartlett.** WRL Technical Note TN-46, March 1995.

''I/O Component Characterization for I/O Cache Designs.'' **Kathy J. Richardson.** WRL Technical Note TN-47, April 1995.

''Attribute caches.'' **Kathy J. Richardson, Michael J. Flynn.** WRL Technical Note TN-48, April 1995.

''Operating Systems Support for Busy Internet Servers.'' **Jeffrey C. Mogul.** WRL Technical Note TN-49, May 1995.

''The Predictability of Libraries.'' **Brad Calder, Dirk Grunwald, Amitabh Srivastava.** WRL Technical Note TN-50, July 1995.

WRL Research Reports and Technical Notes are available on the World Wide Web, from `http://www.research.digital.com/wrl/techreports/index.html`.