

AD-A135 201

THE DISTRIBUTED V KERNEL AND ITS PERFORMANCE FOR
DISKLESS WORKSTATIONS; STANFORD UNIV CA DEPT OF
COMPUTER SCIENCE D R CHERTON ET AL JUL 83
STAN CS 83 973 MDA903 80 C 0102

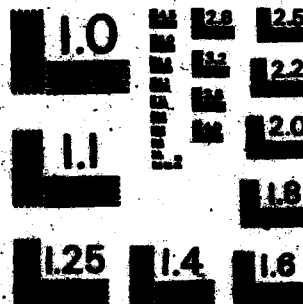
1/1

UNCLASSIFIED

1/6 17/2 NI



END
1/1
1/1



MICROCOPY RESOLUTION TEST CHART
NATIONAL BUREAU OF STANDARDS-1963-A

July 1983

Report No. STAN-CS-83-973
Also numbered CSL 246

②

The Distributed V Kernel and Its Performance for Diskless Workstations

by

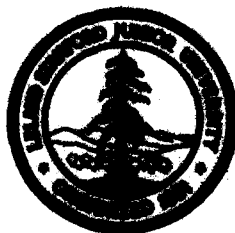
David R. Cheriton and Willy Zwanevool

Contract NDA-903-80-0-0103

Department of Computer Science

Stanford University
Stanford, CA 94305

APPROVED FOR PUBLIC RELEASE
DISTRIBUTION UNLIMITED



DTIC FILE COPY

DTIC
ELECTE
DEC 1 1983
S A D

88 : 10 04 015

AD-A185-201

The Distributed V Kernel and Its Performance for Diskless Workstations

David R. Cheriton and Willy Swenepoel

Computer Systems Laboratory
Departments of Computer Science and Electrical Engineering
Stanford University

The author

Abstract

The distributed V kernel is a microkernel-based system that provides efficient local and network transparent communication. It is primarily being used to do experiments of diskless workstations connected by a high-speed local network to a set of file servers. We describe a performance evaluation of the kernel, with particular emphasis on the cost of network file access. Our results show that over a local network:

1. Diskless workstations can access remote files with minimal performance penalty.
2. The V storage facility can be used to access remote files at comparable cost to a local disk-based storage facility.

It is concluded

that it is feasible to build a distributed system with all network communication using the V storage facility even when most of the network nodes have no permanent storage.

1. Introduction

The distributed V kernel is a microkernel-based system that provides efficient local and network transparent communication. The kernel consists of a set of processes that run on a set of nodes. It is designed to be used in a distributed system where nodes are connected by a high-speed local network. The kernel is designed to be used in a distributed system where nodes are connected by a high-speed local network. The kernel is designed to be used in a distributed system where nodes are connected by a high-speed local network.

Network interprocess communication is predominantly used for remote file access since most SUN workstations at Stanford are configured without a local disk.

This paper reports our experience with the implementation and use of the V kernel. Of particular interest are the experimental results of our approach, namely:

1. The use of diskless workstations with all secondary storage provided by remote file servers.
2. The use of a general purpose network transparent communication facility for support of special-purpose file access protocols and, in particular, the use of a Thrash-like transparent communication mechanism.

The most experimental approach is to configure workstations with a small local disk, using server-based file systems for primary storage. Diskless workstations, however, have a number of advantages including:

1. Lower hardware cost per workstation.
2. Greater responsiveness and simplicity of code with shared file servers.
3. Little or no memory or processing overhead on the workstations for file system and disk handling.
4. Tight coupling with application, reducing the distribution of files.

The major disadvantage is the overhead of processing all file access requests on the server. Our approach involves using the V storage facility to access remote files. The V storage facility is a microkernel-based system that provides efficient local and network transparent communication. It is designed to be used in a distributed system where nodes are connected by a high-speed local network. The V storage facility is designed to be used in a distributed system where nodes are connected by a high-speed local network.

well as relatively easy to use, does not support application-level use of streaming.

These potential problems prompted a performance evaluation of our methods, with particular emphasis on the efficiency of file access. This emphasis on file access distinguishes our work from similar studies [10, 13]. The results of our study strongly support the idea of building a distributed system using distinct workstations connected by a high-speed local network to one or more file servers. Furthermore, we show that stream file access using the V kernel IPC facility is only slightly more expensive than a lower bound imposed by the basic cost of network communication. From this we conclude that relatively little improvement in performance can be achieved using protocols further specialized to file access.

2. V Kernel Interprocess Communication

The basic model provided by the V kernel is that of many small processes communicating by messages. A process is identified by a 32-bit globally unique process identifier or *pid*. Communication between processes is provided in the form of short fixed-length messages, each with an associated reply message, plus a data transfer operation for moving larger amounts of data between processes. In particular, all messages are a fixed 32 bytes in length.

The common communication scenario is as follows: A client process executes a *Send* to a server process which then completes execution of a *Receive* to deliver the message and eventually executes a *Reply* to respond with a reply message back to the client. We refer to this sequence as a *message exchange*. The receiver may execute one of *MoveTo* or *MoveFrom* data transfer operations between the time the message is received and the time the reply message is sent.

The following sections describe the primitives relevant to this paper. The *kernel* refers to software in the V kernel named [5] for a complete description of the kernel facilities.

2.1. Primitives

Send(*message*, *pid*)

Send a message to the process identified by *pid*. The message is a 32-byte fixed-length message. The *message* parameter is a pointer to the message buffer. The *pid* parameter is a 32-bit process identifier. The *Send* operation is a blocking operation. The process calling *Send* will not return until the message has been received by the process identified by *pid*. The *Send* operation is a blocking operation. The process calling *Send* will not return until the message has been received by the process identified by *pid*.

Receive

(*pid*, *scope*) = *Receive*(*message*, *scope*, *reply*, *scope*)
Receive a message from the process identified by *pid*. The message is a 32-byte fixed-length message. The *message* parameter is a pointer to the message buffer. The *scope* parameter is a 32-bit process identifier. The *reply* parameter is a pointer to the reply buffer. The *scope* parameter is a 32-bit process identifier. The *Receive* operation is a blocking operation. The process calling *Receive* will not return until the message has been received from the process identified by *pid*.

Reply(*message*, *pid*)

Send a 32-byte reply message to the process identified by *pid*. The message is a 32-byte fixed-length message. The *message* parameter is a pointer to the message buffer. The *pid* parameter is a 32-bit process identifier. The *Reply* operation is a blocking operation. The process calling *Reply* will not return until the message has been received by the process identified by *pid*.

ReplyWithSegment(*message*, *pid*, *segment*, *scope*, *scope*)

Send a reply message to the process identified by *pid* and transfer the data segment specified by *segment* and *scope* to the process identified by *scope*. The *message* parameter is a pointer to the message buffer. The *pid* parameter is a 32-bit process identifier. The *segment* parameter is a pointer to the data segment. The *scope* parameter is a 32-bit process identifier. The *ReplyWithSegment* operation is a blocking operation. The process calling *ReplyWithSegment* will not return until the message has been received by the process identified by *pid*.

MoveFrom(*segment*, *dest*, *src*, *scope*)

Copy data from the segment starting at *src* in the address space of *scope* to the segment starting at *dest* in the address space of *scope*. The *segment* parameter is a pointer to the data segment. The *dest* parameter is a pointer to the destination address. The *src* parameter is a pointer to the source address. The *scope* parameter is a 32-bit process identifier. The *MoveFrom* operation is a blocking operation. The process calling *MoveFrom* will not return until the data has been copied.

MoveTo(*segment*, *dest*, *src*, *scope*)

Copy data from the segment starting at *src* in the address space of *scope* to the segment starting at *dest* in the address space of *scope*. The *segment* parameter is a pointer to the data segment. The *dest* parameter is a pointer to the destination address. The *src* parameter is a pointer to the source address. The *scope* parameter is a 32-bit process identifier. The *MoveTo* operation is a blocking operation. The process calling *MoveTo* will not return until the data has been copied.

GetFileSegment(*file*, *scope*)

Return a pointer to the segment identified by *file* in the address space of *scope*. The *file* parameter is a pointer to the file descriptor. The *scope* parameter is a 32-bit process identifier. The *GetFileSegment* operation is a blocking operation. The process calling *GetFileSegment* will not return until the segment has been identified.

pid = *GetFileSegment*(*file*, *scope*)

Return the process identifier associated with *file* in the address space of *scope*. The *file* parameter is a pointer to the file descriptor. The *scope* parameter is a 32-bit process identifier. The *pid* operation is a blocking operation. The process calling *pid* will not return until the process identifier has been identified.

2.2. Operations

The V kernel provides a set of operations for moving data between processes. These operations are: *Send*, *Receive*, *Reply*, *ReplyWithSegment*, *MoveFrom*, *MoveTo*, *GetFileSegment*, and *pid*. These operations are used to move data between processes. The *Send* operation is used to send a message to a process. The *Receive* operation is used to receive a message from a process. The *Reply* operation is used to send a reply message to a process. The *ReplyWithSegment* operation is used to send a reply message to a process and transfer data to the process. The *MoveFrom* operation is used to copy data from one segment to another. The *MoveTo* operation is used to copy data from one segment to another. The *GetFileSegment* operation is used to get a pointer to a segment. The *pid* operation is used to get the process identifier of a process.

The V kernel provides a set of operations for moving data between processes. These operations are: *Send*, *Receive*, *Reply*, *ReplyWithSegment*, *MoveFrom*, *MoveTo*, *GetFileSegment*, and *pid*. These operations are used to move data between processes. The *Send* operation is used to send a message to a process. The *Receive* operation is used to receive a message from a process. The *Reply* operation is used to send a reply message to a process. The *ReplyWithSegment* operation is used to send a reply message to a process and transfer data to the process. The *MoveFrom* operation is used to copy data from one segment to another. The *MoveTo* operation is used to copy data from one segment to another. The *GetFileSegment* operation is used to get a pointer to a segment. The *pid* operation is used to get the process identifier of a process.

The V kernel provides a set of operations for moving data between processes. These operations are: *Send*, *Receive*, *Reply*, *ReplyWithSegment*, *MoveFrom*, *MoveTo*, *GetFileSegment*, and *pid*. These operations are used to move data between processes. The *Send* operation is used to send a message to a process. The *Receive* operation is used to receive a message from a process. The *Reply* operation is used to send a reply message to a process. The *ReplyWithSegment* operation is used to send a reply message to a process and transfer data to the process. The *MoveFrom* operation is used to copy data from one segment to another. The *MoveTo* operation is used to copy data from one segment to another. The *GetFileSegment* operation is used to get a pointer to a segment. The *pid* operation is used to get the process identifier of a process.

must be allocated in the kernel and large amounts of data are transferred directly between user address spaces without user copies. Moreover, by virtue of the sparsity of the communication, the kernel's message buffers can be naturally allocated. As exemplified in Thrash, these factors make for a small, efficient kernel.

The V message primitives appear ill-suited in several ways for a network environment, at least on first observation. The short, fixed-length messages appear to make inefficient use of large packet sizes typically available on local networks. The synchronous nature of the message sending would seem to interfere with the true parallelism possible between separate workstations. And the economics of message buffering afforded by these restrictions in a single machine implementation are less evident in a distributed environment. Finally, the separate data transfer operations *MoveTo* and *MoveFrom* appear only to increase the number of remote data transfer operations that must be implemented in the distributed case.

However, our experience has been that the V message primitives are easily and efficiently implemented over a local network. Moreover, we have found that the messages of the primitives facilitated an efficient distributed implementation. The only major departure from Thrash was the explicit specification of segments in messages and the addition of the primitives *ReceiveWithSegment* and *ReplyWithSegment*. This extension was done for efficient page-level file access although it has proven useful under more general circumstances, eg. in pushing disaster relief notices to remote servers.

3. Implementation Issues

A foremost concern in the implementation of the kernel has been efficiency. Before describing some of the implementation details of the individual primitives, we first present aspects of the implementation that are common to the different operations of the kernel.

1. Remote operations are implemented directly in the kernel instead of through a process-level message layer. When the kernel receives a request from a remote process, it immediately copies the request into a kernel buffer. The kernel then copies the request into a kernel buffer and immediately sends the request to the remote process. This is done by the kernel's *Send* primitive. The kernel then copies the request into a kernel buffer and immediately sends the request to the remote process. This is done by the kernel's *Send* primitive. The kernel then copies the request into a kernel buffer and immediately sends the request to the remote process. This is done by the kernel's *Send* primitive.

interest functionality and local net performance, we have chosen not to include the dominant (but not) operation with any more overhead than is strictly necessary.

1. The synchronous request-response nature of a reply associated with each message is exploited to build reliable message transmission directly on an unreliable transport service, i.e. without using an extra layer (and extra packets) to implement reliable transport. The reply message serves as an acknowledgment as well as carrying the reply message data.
4. The mapping from process id to process location is aided by encoding a host specification in the process identifier. The kernel can then determine quickly whether a process is either local or remote, and in the latter case on which machine it resides.
5. There are no per-packet acknowledgments for large data transfers (as in *MoveTo* and *MoveFrom*). There is only a single acknowledgment when the transfer is complete.
6. File page-level transfer requires the optional number of packets (i.e. two) because of the ability to append short segments to messages using *ReceiveWithSegment* and *ReplyWithSegment*.

The following sections look at particular aspects of the implementation in greater detail.

3.1. Process Naming

V uses a global (but) naming space for specifying processes, in contrast to the local port naming used in DESOS [1] and Aegis [2]. Process identifiers are unique within the context of a host network. On the SUN workstation, it is natural for the V kernel to use 32-bit process identifiers. The high-order 16 bits of the process identifier serve as a logical host identifier which the low-order 16 bits are used as a locally unique identifier.

In the current 32-bit identifier implementation, the top 8 bits of the logical host identifier are the physical network address of the workstation, making the process identifier an external address mapping device. In the 16-bit implementation, a 16-bit range high-order to network address. When there is no such entry for the specified logical host, the message is dropped. Now "logical-destination address" correspondence can be determined from message received. However, with such small or large range of logical to physical host identifier being used.

The use of an explicit host field in the process identifier allows distributed processing of remote process identifier requests. The kernel can then, in addition to using the high-order 16 bits of the process identifier to determine the logical host, use the low-order 16 bits to determine the physical host. This allows the kernel to determine the physical host of a process without having to maintain a table of physical host identifiers. In addition, the V kernel can use the high-order 16 bits of the process identifier to determine the logical host of a process without having to maintain a table of logical host identifiers.

GoM uses network browsing to determine the mapping of a logical process identifier to a physical identifier if the mapping is not known to the local kernel. Any kernel handling the mapping are reported to the browsing engine. The addition of local and remote scrapes was required to discriminate between server processes that serve only a single workstation and those that serve the network.

When a process identifier is specified to `Send` with a logical host identifier different from that of the local machine, the local pid validation test fails and `Send` calls `NonLocalSend` which handles transmission of the message over the network.

THE FOLLOWING INFORMATION IS FOR YOUR INFORMATION ONLY. IT IS NOT TO BE USED FOR ANY OTHER PURPOSE. IT IS NOT TO BE REPRODUCED OR TRANSMITTED IN ANY FORM OR BY ANY MEANS, ELECTRONIC OR MECHANICAL, INCLUDING PHOTOCOPYING, RECORDING, OR BY ANY INFORMATION STORAGE AND RETRIEVAL SYSTEM, WITHOUT THE WRITTEN PERMISSION OF THE NATIONAL ARCHIVES.

As in the local case, major connections arise with *AddressTo* and *AddressFrom* because, by their definition, there is no need for grouping or building of the data in the kernel. The V kernel sends the data directly from the source address space into the network interface, and directly from the network interface to the destination address space.

The message and data transfer primitives provide efficient communication of small amounts and large amounts of data, less than 32 bytes or several kilobytes of network packets. However, perhaps the unique system feature is an intermediate amount of data that is not efficiently handled by the Thrift primitives when communicating over a network.

[illegible]

Activities of the Interviewer

1. The advantages of the Thoth IPC model with the enhanced performance of a WFE-type file access protocol.
2. Compatibility with efficient local operation. For business segments may be accessed directly if it is the most efficient space or if the recipient process operates a DMA device.
3. Use of Access/STASegments and Read/STASegments is automatic and transparent to application programs. (See Table 1)

An expected objection to our solution is the apparent redundancy and self-redundancy of the message problem. We find the redundancy may be overcompensating for the persistent redundancy of communication and synchronization between client and server processes. Also, it is not unexpected that there is some overlap between efficiency and fairness. Given the requirements currently imposed upon system services through such systems that provide a pseudonym interface to the message problem, it is not inappropriate to make these comparisons as the problem itself is efficient.

We now turn to the discussion of the performance evaluation of the board. We first define the term *board quality* as a reasonably lower bound to the state of board composition. Subsequently we discuss the efficiency of the board, both in terms of various policies and decisions.

[illegible][illegible]

1. The cost of private operations versus the cost of the corresponding public provision.
2. The cost of the extra costs of private versus public provision.

[illegible]

network. The network gateway is a function of the processor, the network, the network interface and the number of links transferred. It is the critical gate point for interpreting the network, however, test software modules that could otherwise transmit the data by proxy, packet. The network gateway is designed by connecting the data to transmit a signal from the main memory of one computer to the main memory of another and vice versa and dividing the cost flow for the equivalent by 2. The equivalent is represented a large number of times for statistical purposes. The modules are implemented at the data link layer using the language tool so that no program or program-switching overhead appears in the code. The assumption of operation throughout and low network utilization are good characteristics of most test network environments.

However, primarily providing a more realistic minimum caloric intake than the transfer than that suggested by the physical amount, again because it includes the processes and transport interface steps. For instance, a 10 kWh thermal can mean 1000 kWh from one combustion to another in 100 milliseconds. However, the time to transport the product to the interface at the sending end and the time to transfer the product out of the interface at the receiving end are comparable to the time for combustion. Thus, the time for the transfer from the point of view of the compensating surface transfer is at least one or three times as long as the transfer by the physical transfer one.

Measurement of network quality was made using the experimental 2 dB network. In all measurements, the network was normally idle due to the considerable time at which measurements had to be made. Table 4-1 lists our measurements of the 2 dB network quality for the SMC configuration using the 2 and 40 dBm transmitters with those given in Table 4-2. The network measurements prove that the time for the data to be transmitted based on the physical bit rate of the medium, namely 2.54 Mbit/sec, is adequate.

CONFIDENTIAL

Speed	Network Time	Network Capacity	
		1990s	2000s
0.5	174	0.50	0.50
1.0	348	1.00	1.00
1.5	522	1.50	1.50
2.0	696	2.00	2.00
2.5	870	2.50	2.50

SECRET

The difference between the network data, computed in the network data file, and the estimated network quality data is accounted for primarily by the parameter that is present and controls the quality and thus relative the quality of the other end. For instance, with a 100% type packet and an 80% type packet, the copy time from memory to the network interface and vice versa is roughly 1.50 milliseconds in each direction. This, of the total 6.95 milliseconds, 1.50 is copy time, 5.45 is network transmission time and 3 is (approximately) network congestion latency. If we transfer a 100% packet with double latency, the transmission time is less than 2 milliseconds while the copy times remain the same, reducing the parameter that 75 percent of the cost is the network latency. The importance of the parameter speed is also illustrated by the difference in network quality for the two processors measured in Table 4-1.

With our interface, the parameter is required to copy the packet into the interface for transmission and this of the interface as reception (with the interface providing bandwidth as buffer buffering). From the copy time given above, we might expect that a DMA interface would significantly improve performance. However, we would predict that a DMA interface would not result in better latency performance for the network. This is because the network interface is a copy of the packet from the network interface, allowing it to be placed back at the packet data immediately in the DMA interface. With a DMA interface, the copy would be made directly to the packet data from DMA's last main memory. Similarly, as transmission, the network interface is a copy of the packet from the network interface. With DMA interface, copies the packet to the parameter in the computer area of memory, keeping the main data in the network copy operation. Finally, there is no network interface in the network interface, a network DMA interface. The network interface moves data faster than a 100% DMA interface. This parameter is used in the DMA interface. In general, the main benefit of a main network interface appears to be in allowing the main processor rather than spending its operation for main use of network operations.

5. Kernel Performance

Our first kernel performance study is the question of network interface performance. The main performance is provided in terms of the data in Table 4-1. The network interface is a copy of the packet from the network interface, allowing it to be placed back at the packet data immediately in the DMA interface. With a DMA interface, the copy would be made directly to the packet data from DMA's last main memory. Similarly, as transmission, the network interface is a copy of the packet from the network interface. With DMA interface, copies the packet to the parameter in the computer area of memory, keeping the main data in the network copy operation. Finally, there is no network interface in the network interface, a network DMA interface. The network interface moves data faster than a 100% DMA interface. This parameter is used in the DMA interface. In general, the main benefit of a main network interface appears to be in allowing the main processor rather than spending its operation for main use of network operations.

Measurement of primary interface was done using a low-priority "background" process on each workstation that repeatedly copies a packet to or from the loop. All other processor activities follow the processor schedule in the packet. Thus, the packet that used per operation as a workstation is the elapsed time since the processor that started in the "background" process divided by 24, the number of operations executed.

Using 1000 words per operation and data which give or return 10 milliseconds, our measurements should be accurate to about 20 milliseconds except for the effect of variation in network load.

5.1. Kernel Measurements

Table 5-1 gives the results of our measurements of average network load data similar with the kernel coming on workstation using a 100% processor and measured by a 3.40 seconds. This data shows that a kernel network operation, including the time to copy the packet into the network interface, is about 10 milliseconds. This number is about 10 milliseconds and shows the effect of the network interface on the network interface. The network interface is a copy of the packet from the network interface, allowing it to be placed back at the packet data immediately in the DMA interface. With a DMA interface, the copy would be made directly to the packet data from DMA's last main memory. Similarly, as transmission, the network interface is a copy of the packet from the network interface. With DMA interface, copies the packet to the parameter in the computer area of memory, keeping the main data in the network copy operation. Finally, there is no network interface in the network interface, a network DMA interface. The network interface moves data faster than a 100% DMA interface. This parameter is used in the DMA interface. In general, the main benefit of a main network interface appears to be in allowing the main processor rather than spending its operation for main use of network operations.

A second test of network interface was done using a low-priority "background" process on each workstation that repeatedly copies a packet to or from the loop. All other processor activities follow the processor schedule in the packet. Thus, the packet that used per operation as a workstation is the elapsed time since the processor that started in the "background" process divided by 24, the number of operations executed.

Kernel Performance

Kernel Operation	Elapsed Time			Network Penalty	Processor Time	
	Local	Remote	Difference		Client	Server
GetTime	0.07	-	-	-	0.07	-
Send-Receive-Reply	1.00	1.18	2.18	1.00	1.79	2.30
MoveFrom: 1024 bytes	1.26	9.83	7.77	8.15	3.76	5.69
MoveTo: 1024 bytes	1.26	9.85	7.79	8.15	3.59	5.87

Table 5-1: Kernel Performance: 3 MB Ethernet and 80 MHz Processor (times in milliseconds)

Kernel Operation	Elapsed Time			Network Penalty	Processor Time	
	Local	Remote	Difference		Client	Server
GetTime	0.06	-	-	-	0.06	-
Send-Receive-Reply	0.77	2.54	1.77	1.30	1.44	1.79
MoveFrom: 1024 bytes	0.95	8.60	7.65	6.77	3.32	4.78
MoveTo: 1024 bytes	0.95	8.60	7.65	6.77	3.17	4.95

Table 5-2: Kernel Performance: 3 MB Ethernet and 10 MHz Processor (times in milliseconds)

view interprocess communication as transparent across machines when the speed ratio is so large. However, an alternative interpretation is to recognize that the remote operation adds a delay of less than 2 milliseconds, and that in many cases this delay is insignificant relative to the time necessary to generate a request in the server. Furthermore, the waiting at client workstation processor is tiny with the server load for only 1.04 milliseconds out of the total 2.14 milliseconds time taken by the 80 MHz processor). Then, one can afford the processor on one machine to, for instance, moving a server process to another machine if its request processing generally requires more than 0.07 milliseconds of processor time, i.e. the difference between the local time for Send-Receive-Reply and the local penalty time for the delay operation.

5.4. Small-Program Traffic

The situation so far has been for a single pair of processes communicating over the network. Operating systems and network configurations would be using the network infrastructure to communicate with other programs. Some workstations to implement distributed systems management and control the network and the operation is communication with other workstations.

A pair of workstations connected by a network can be used to implement a distributed system. The network infrastructure would be used to communicate with other programs. Some workstations to implement distributed systems management and control the network and the operation is communication with other workstations.

network in this fashion. Unfortunately, our measurements of this scenario turned up a hardware bug in our 3 MB Ethernet interface, a bug which caused many collisions to go undetected and show up as corrupted packets. Thus, the response time for the 8 MHz processor workstation in this case is 3.4 milliseconds. The increase in time from 1.18 milliseconds is accounted for almost entirely by the detected and retransmitted data (usually one per 200 packets) from the hardware bug. With standard network interfaces, we believe that the network can support any reasonable level of message communication without significant performance degradation.

A more critical resource is processor time. This is especially true for workstations which are servers that need to be the focus of a distributed network of message traffic. For instance, just based on workstations that a workstation is loaded to a point about 50% message exchanges per second, independent of the number of clients. The number is substantially lower for the client workstations, particularly when a remote client for the server processing is involved. The workstations are connected to the network.

The network infrastructure would be used to communicate with other programs. Some workstations to implement distributed systems management and control the network and the operation is communication with other workstations.

We first describe the performance of random page-level file access.

Table 6-1 list the times for reading or writing a 512 byte block between two processes both local and remote using the 10 MHz processor. The times do not include time to fetch the data from disk but do indicate expected performance when data is buffered in memory. A page send involves the sequence of kernel operations: *Send-Receive-ReplyWithSegment*. A page write is *Send-ReceiveWithSegment-Reply*.

These measurements indicate the performance when file reading and writing use explicit segment specification in the *name* and *ReceiveWhatSegment* and *ReplyWhatSegment*. However, a file write can also be performed in a more basic Thrift-like way using the *Send-Receive-MoveFrom-Reply* sequence. For a 512 byte write, this costs 8.1 milliseconds; file reading is similar using *MoveTo*. Thus, the segment mechanism saves 3.5 milliseconds on every page read and write operation, justifying this extension to the machine primitives.

Operation	Elapsed Time			Network Penalty	Processor Time	
	Local	Remote	FileShare		Client	Server
page read	1.31	3.96	4.25	3.89	2.50	3.28
page write	1.31	5.00	4.29	3.89	2.58	3.32

The columns are to be interpreted according to the explanation given for similarly labeled columns of Tables 3-1 and 3-2. Thus, the data to read at the top of a page using these procedures is approximately 1.5 milliseconds more than the current latency for these conditions.

The average savings for a total tape system, compared to a tape drive in a computer system, is about 10 percent. Savings is also being realized in the time taken to load and unload tape. The time required for disk loading has come with mounting cost. If the average disk costs \$100 for a 5 1/4 inch 5.25 megabyte disk then the average load disk costs \$100 per hour. That is up from \$50 per hour in 1980. The cost of tape is still less than the cost of disk.

Segmented file access is the predominant pattern for file activity in most systems. Efficient file systems exploit this behavior to reduce the effect of disk latency by prefetching file pages (read-ahead) and asynchronously flushing modified pages (write-behind). File access and file transfer protocols typically implement streaming to reduce the effect of network latency on performance of the system.

Two factors suggest that increasing muscle development is a local phenomenon. First, most animals have a low degree of a concentration of fast-twitch speed and low-twitch "stomch" strength. Second, 100 to 150% is highly significant, significant increases occur in the extreme musculature. Further evidence is given by research findings. The pattern of muscle development is very consistent in groups of older mice, but very slow and temporary of rate. However, building muscle is a very slow rate in the immature, but when the young animals mature to adult, the rate of muscle development is very rapid.

1. The first step in the process is to identify the problem or issue that needs to be addressed. This involves gathering information and understanding the context of the problem.

1. The first step in the process is to identify the problem. This involves gathering information about the situation and understanding the needs of the stakeholders involved.

1. Subject: [REDACTED]
 2. Reference: [REDACTED]
 3. Summary: [REDACTED]
 4. Remarks: [REDACTED]
 5. Signature: [REDACTED]
 6. Date: [REDACTED]
 7. Initials: [REDACTED]
 8. Comments: [REDACTED]
 9. Remarks: [REDACTED]
 10. Signature: [REDACTED]
 11. Date: [REDACTED]
 12. Initials: [REDACTED]
 13. Comments: [REDACTED]
 14. Remarks: [REDACTED]
 15. Signature: [REDACTED]
 16. Date: [REDACTED]
 17. Initials: [REDACTED]
 18. Comments: [REDACTED]
 19. Remarks: [REDACTED]
 20. Signature: [REDACTED]
 21. Date: [REDACTED]
 22. Initials: [REDACTED]
 23. Comments: [REDACTED]
 24. Remarks: [REDACTED]
 25. Signature: [REDACTED]
 26. Date: [REDACTED]
 27. Initials: [REDACTED]
 28. Comments: [REDACTED]
 29. Remarks: [REDACTED]
 30. Signature: [REDACTED]
 31. Date: [REDACTED]
 32. Initials: [REDACTED]
 33. Comments: [REDACTED]
 34. Remarks: [REDACTED]
 35. Signature: [REDACTED]
 36. Date: [REDACTED]
 37. Initials: [REDACTED]
 38. Comments: [REDACTED]
 39. Remarks: [REDACTED]
 40. Signature: [REDACTED]
 41. Date: [REDACTED]
 42. Initials: [REDACTED]
 43. Comments: [REDACTED]
 44. Remarks: [REDACTED]
 45. Signature: [REDACTED]
 46. Date: [REDACTED]
 47. Initials: [REDACTED]
 48. Comments: [REDACTED]
 49. Remarks: [REDACTED]
 50. Signature: [REDACTED]
 51. Date: [REDACTED]
 52. Initials: [REDACTED]
 53. Comments: [REDACTED]
 54. Remarks: [REDACTED]
 55. Signature: [REDACTED]
 56. Date: [REDACTED]
 57. Initials: [REDACTED]
 58. Comments: [REDACTED]
 59. Remarks: [REDACTED]
 60. Signature: [REDACTED]
 61. Date: [REDACTED]
 62. Initials: [REDACTED]
 63. Comments: [REDACTED]
 64. Remarks: [REDACTED]
 65. Signature: [REDACTED]
 66. Date: [REDACTED]
 67. Initials: [REDACTED]
 68. Comments: [REDACTED]
 69. Remarks: [REDACTED]
 70. Signature: [REDACTED]
 71. Date: [REDACTED]
 72. Initials: [REDACTED]
 73. Comments: [REDACTED]
 74. Remarks: [REDACTED]
 75. Signature: [REDACTED]
 76. Date: [REDACTED]
 77. Initials: [REDACTED]
 78. Comments: [REDACTED]
 79. Remarks: [REDACTED]
 80. Signature: [REDACTED]
 81. Date: [REDACTED]
 82. Initials: [REDACTED]
 83. Comments: [REDACTED]
 84. Remarks: [REDACTED]
 85. Signature: [REDACTED]
 86. Date: [REDACTED]
 87. Initials: [REDACTED]
 88. Comments: [REDACTED]
 89. Remarks: [REDACTED]
 90. Signature: [REDACTED]
 91. Date: [REDACTED]
 92. Initials: [REDACTED]
 93. Comments: [REDACTED]
 94. Remarks: [REDACTED]
 95. Signature: [REDACTED]
 96. Date: [REDACTED]
 97. Initials: [REDACTED]
 98. Comments: [REDACTED]
 99. Remarks: [REDACTED]
 100. Signature: [REDACTED]
 101. Date: [REDACTED]
 102. Initials: [REDACTED]
 103. Comments: [REDACTED]
 104. Remarks: [REDACTED]
 105. Signature: [REDACTED]
 106. Date: [REDACTED]
 107. Initials: [REDACTED]
 108. Comments: [REDACTED]
 109. Remarks: [REDACTED]
 110. Signature: [REDACTED]
 111. Date: [REDACTED]
 112. Initials: [REDACTED]
 113. Comments: [REDACTED]
 114. Remarks: [REDACTED]
 115. Signature: [REDACTED]
 116. Date: [REDACTED]
 117. Initials: [REDACTED]
 118. Comments: [REDACTED]
 119. Remarks: [REDACTED]
 120. Signature: [REDACTED]
 121. Date: [REDACTED]
 122. Initials: [REDACTED]
 123. Comments: [REDACTED]
 124. Remarks: [REDACTED]
 125. Signature: [REDACTED]
 126. Date: [REDACTED]
 127. Initials: [REDACTED]
 128. Comments: [REDACTED]
 129. Remarks: [REDACTED]
 130. Signature: [REDACTED]
 131. Date: [REDACTED]
 132. Initials: [REDACTED]
 133. Comments: [REDACTED]
 134. Remarks: [REDACTED]
 135. Signature: [REDACTED]
 136. Date: [REDACTED]
 137. Initials: [REDACTED]
 138. Comments: [REDACTED]
 139. Remarks: [REDACTED]
 140. Signature: [REDACTED]
 141. Date: [REDACTED]
 142. Initials: [REDACTED]
 143. Comments: [REDACTED]
 144. Remarks: [REDACTED]
 145. Signature: [REDACTED]
 146. Date: [REDACTED]
 147. Initials: [REDACTED]
 148. Comments: [REDACTED]
 149. Remarks: [REDACTED]
 150. Signature: [REDACTED]
 151. Date: [REDACTED]
 152. Initials: [REDACTED]
 153. Comments: [REDACTED]
 154. Remarks: [REDACTED]
 155. Signature: [REDACTED]
 156. Date: [REDACTED]
 157. Initials: [REDACTED]
 158. Comments: [REDACTED]
 159. Remarks: [REDACTED]
 160. Signature: [REDACTED]
 161. Date: [REDACTED]
 162. Initials: [REDACTED]
 163. Comments: [REDACTED]
 164. Remarks: [REDACTED]
 165. Signature: [REDACTED]
 166. Date: [REDACTED]
 167. Initials: [REDACTED]
 168. Comments: [REDACTED]
 169. Remarks: [REDACTED]
 170. Signature: [REDACTED]
 171. Date: [REDACTED]
 172. Initials: [REDACTED]
 173. Comments: [REDACTED]
 174. Remarks: [REDACTED]
 175. Signature: [REDACTED]
 176. Date: [REDACTED]
 177. Initials: [REDACTED]
 178. Comments: [REDACTED]
 179. Remarks: [REDACTED]
 180. Signature: [REDACTED]
 181. Date: [REDACTED]
 182. Initials: [REDACTED]
 183. Comments: [REDACTED]
 184. Remarks: [REDACTED]
 185. Signature: [REDACTED]
 186. Date: [REDACTED]
 187. Initials: [REDACTED]
 188. Comments: [REDACTED]
 189. Remarks: [REDACTED]
 190. Signature: [REDACTED]
 191. Date: [REDACTED]

SECRET

100-443887-100

1. The first step in the process is to identify the problem or issue that needs to be addressed. This involves gathering information and understanding the context of the problem.



many current time-sharing systems do, such program loading could achieve the same performance given in the table, independent of disk speed. Thus, we argue that *Advent* and *More/From* with large transfer units provide an efficient program loading mechanism that is as fast as can be achieved with the given hardware.

7. File Server Issues

File server performance is a critical issue for desktop workstations. Unfortunately, we do not yet have experience with a V kernel-based file server. Thus, this section describes what we believe are the key issues and estimates performance without providing conclusive data. In general, we view the processor as the key resource to consider in file server performance because, as argued earlier, the network bandwidth is plentiful and disk scheduling and buffering issues are identical to those encountered in conventional multi-user systems.

The number of workstations a file server can support can be estimated from processor requirements. If we estimate page read or write processing overhead as roughly 15 milliseconds for file system processing (from LOCUS) plus 15 milliseconds for kernel operation (from Table 6-1), a page request costs about 7 milliseconds of processor time. Program loading appears to cost about 300 milliseconds for an average 64 kilobyte program. Estimating that 50 percent of the file requests are page requests, the average request costs 26 milliseconds. Thus, a file server based on the SUN workstation processor could support about 38 file requests a second. From this we estimate that one file server can serve about 20 workstations adequately, but 25 or more active workstations would lead to excessive delays. However, a desktop workstation system can easily be extended to handle more workstations by adding more file server machines along the network would not seem to be a bottleneck for less than 200 workstations.

For some programs, it is advantageous to view of the kernel processor requirements to estimate the processor for the server, rather than to find the processor has a minimum and subsequently add more page requests to the processor, causing for a slow start and doing a lot of the work with entering out in the data requests if they require substantial interaction with the user.

On this basis, a file server should have a general purpose processor capable of the ability to identify source code programs. The cost for the processor is \$1000 to \$1500. A further argument for using a general purpose processor is that the cost of a general purpose processor is less than the cost of a specialized processor. The cost of a specialized processor is \$1000 to \$1500. A further argument for using a general purpose processor is that the cost of a general purpose processor is less than the cost of a specialized processor. The cost of a specialized processor is \$1000 to \$1500. A further argument for using a general purpose processor is that the cost of a general purpose processor is less than the cost of a specialized processor. The cost of a specialized processor is \$1000 to \$1500.

program, i.e. it is transparent except for performance.

8. Measurements with the 10 Mb Ethernet

Our limited access to a 10 Mb Ethernet has precluded basing our measurements on this standard local network. However, some preliminary figures using the 10 Mb Ethernet indicate the effect of using a faster network and slightly faster network interfaces. First, the remote message exchange time is 2.71 milliseconds using an 8 MHz processor, roughly the time for the 10 MHz processor on the 3 Mb network and .5 milliseconds better than the 8 MHz processor on the 3 Mb network. Second, the page read time is 5.72 milliseconds. Finally, the program loading time is much improved, achieving 255 milliseconds for a 64 kilobyte load using 16 Kb transfer units. We have not identified to what degree the improvement is due to the faster network speed versus the difference in the network interface.

9. Related Work

There are a number of previous and concurrent efforts in providing communication mechanisms for distributed systems. For brevity, we compare our work with only a representative sample that characterizes the search for, and evaluation of, reliable models and implementations.

Spencer's remote reference study [13] considered the feasibility of implementing remote load and spare operations over a local network. Nelson's work on remote procedure calls [10] investigates network communication for procedure-based systems analogous to what the V kernel provides for message-based systems. Radich and Robinson's kernel [12] implements a message system with a number of features such as non-blocking message sending that are not provided by the V kernel. Finally, LOCUS [14] implements network communication into a UNIX-like system in the form of a network-based file system.

Our work has followed the general approach of Spencer's and Nelson's work in using a message-based form of communication in place of sending and receiving every program, kernel, or message performance. However, we share Spencer's concern that the system which requires globally shared memory is a distributed system, which is especially the case for the system provided by the processor. The V kernel provides a strong degree of separation between processes and achieves protected provision of services in a multi-user, multi-workstation environment by limiting communication to the kernel PC processor.

Our approach differs from Nelson's primarily in our use of a message-based communication mechanism. This provides a more secure, reliable, and clear of the sharing differences in message. Furthermore, the V kernel provides a system in which to build a robust protection, all although by the addition of message-based and message-based. Under such conditions, the system is more secure and reliable than the system provided by the processor.

versus message-based systems, although it is not clear these differences result in any significant difference in overall performance.

The V kernel performance is roughly comparable to that of the software implementations developed by Spector and Nelson, allowing for the non-trivial differences in operation semantics and host processors. We would hypothesize that V kernel performance could be improved by a factor of 30 using microcode, similar to the improvement observed by Spector and Nelson for their primitives. Unfortunately, neither Spector nor Nelson provides results that afford a comparison with our file access results. In general, their work has concentrated on the speed of the basic mechanism and has not been extended to measure performance in a particular application setting.

In comparison to Accent, the V kernel provides a primitive form of message communication, and benefits accordingly in terms of speed, small code size and ability to run well on an inexpensive machine⁵ without disk or microcode support. For instance, Accent messages require an underlying transport protocol for reliable delivery because there is no client-level reply message associated with every *Send* as in the V kernel. We do not at this time have performance figures for Accent.

LOCUS does not attempt to provide applications with general network interprocess communication but exploits carefully honed problem-oriented protocols for efficient remote file access. It is difficult to compare the two systems from measurements available given the differences in network speeds, processor speeds and measurement techniques. However, from the specific comparisons with LOCUS presented earlier, we would expect overall file access performance for the V kernel to be comparable to LOCUS running on the same machines and network.

However, the memory requirements for the V kernel are about half that of LOCUS compiled for the PDP-11 and probably more like one fourth when LOCUS is compiled for a 32-bit processor like the 68000. Thus, for graphics workstations or process control applications, for instance, the V kernel would be more attractive because of its smaller size, real-time orientation and its provision of general interprocess communication. However, the V kernel does not provide all the functionality of the LOCUS kernel which includes that of the UNIX kernel and more. When required with V, these additional facilities must be provided by server processes executing either on client workstations or network server machines.

10. Conclusions

We conclude that it is feasible to build a distributed system using diskless workstations connected by a high-speed local network to one or more file servers using the V kernel IPC. In particular, the performance study shows that V kernel IPC provides satisfactory

performance despite its generality. Because the performance is so close to the lower bound given by the network penalty, there is relatively little room for improvement on the V IPC for the given hardware regardless of protocol and implementation used.

The efficiency of file access using the V IPC suggests that it can not only replace page-level file access protocols but also file transfer and remote terminal protocols, thereby reducing the number of protocols needed. We claim that V kernel IPC is adequate as a transport level for all our local network communication providing each machine runs the V kernel or at least handles the interkernel protocol. We do, however, see a place for these specific protocols in internetworking situations.

In addition to quantifying the elapsed time for various operations, our study points out the importance of considering processor requirements in the design of distributed systems. More experience and measurement of file server load and workstation file access behavior is required to decide whether file server processing is a significant problem in using diskless workstations.

The V kernel has been in use with the diskless SUN workstations, providing local and remote interprocess communication, since September 1982. It is currently 38 kilobytes including code, data and stack. The major use of the network interprocess communication is for accessing remote files. Our file servers are currently 6 VAX/UNIX systems running a kernel simulator and file server program which provides access to UNIX system services over the Ethernet using interkernel packets. A simple command interpreter program allows programs to be loaded and run on the workstations using these UNIX servers. Our experience with this software to date supports the conclusions of the performance study that we can indeed build our next generation of computing facilities [8] using diskless workstations and the V kernel.

Acknowledgements

We are indebted to all the members of the V research group at Stanford, which at this point includes two faculty members and roughly ten graduate students. In particular, we wish to thank Keith Lauts for his patient comments on a seemingly endless sequence of drafts and Tim Mann for his many contributions to the design and the implementation of the kernel. We would also like to thank the referees whose comments and suggestions helped to enhance the clarity of the paper.

References

1. F. Buxton, J.H. Howard and J.T. Manning. Task Communication in DEDACS. Proceedings of the 8th Symposium on Operating System Principles, ACM, November, 1977, pp. 23-31. Publication (Operating Systems Review 12(3)).
2. A. Bach, F. Buxton, V. Paxon. The SUN Workstation Architecture. Tech. Rep. 82, Computer Science Laboratory, Stanford University, March, 1982.

⁵ In addition to the workstations used for this study, the V kernel runs on a PDP-11.

3. D.R. Cheriton, M.A. Malcolm, I.S. Melen and G.R. Sagar. "Thoth, a Portable Real-time Operating System." *Comm. ACM* 22, 2 (February 1979), 105-115.

4. D.R. Cheriton. Distributed I/O using an Object-based Protocol. Tech. Rept. 81-1, Computer Science, University of British Columbia, 1981.

5. D.R. Cheriton. *The Thoth System: Multi-process Structuring and Portability*. American Elsevier, 1982.

6. D.R. Cheriton, T.P. Mann and W. Zwanevool. V-System: Kernel Manual. Computer Systems Laboratory, Stanford University.

7. Digital Equipment Corporation, Intel Corporation and Xerox Corporation. The Ethernet: A Local Area Network - Data Link Layer and Physical Layer Specifications, Version 1.0.

8. K.A. Lantz, D.R. Cheriton and W.I. Nowicki. Third Generation Graphics for Distributed Systems. Tech. Rept. STAN-CS-82-938, Department of Computer Science, Stanford University, February, 1983. To appear in *ACM Transactions on Graphics*.

9. R.M. Metcalfe and D.R. Boggs. "Ethernet: Distributed Packet Switching for Local Computer Networks." *Comm. ACM* 19, 7 (July 1976), 395-404. Also CSL-75-7, Xerox Palo Alto Research Center, reprinted in CSL-80-2.

10. B.J. Nelson. *Remote Procedure Call*. Ph.D. Th., Carnegie-Mellon U., 1981. published as CMU technical report CMU-CS-81-119.

11. G. Popek, B. Walker, J. Chow, D. Edwards, C. Kline, G. Rudisin, G. Thiel. LOCUS: A Network Transparent, High Reliability Distributed System. Proceedings of the 8th Symposium on Operating Systems Principles, ACM, December, 1981, pp. 169-177.

12. R. Rathid and G. Robertson. Accent: A Communication Oriented Network Operating System Kernel. Proceedings of the 8th Symposium on Operating Systems Principles, ACM, December, 1981, pp. 64-75.

13. A. Spitzer. "Performing Remote Operations Efficiently on a Local Computer Network." *Comm. ACM* 25, 4 (April 1982), 246-250.

14. D. Swinehart, G. McDowell and D. Boggs. WFS: A Simple Shared File System for a Distributed Environment. Proceedings of the 7th Symposium on Operating Systems Principles, ACM, December, 1979, pp. 9-23.



Accession For	
DTIC GRAFI	<input checked="" type="checkbox"/>
DTIC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	<input type="checkbox"/>
Distribution/Availability	
Availability Codes	
Avail and/or	Special

DATE
FILMED
8