

# Hierarchical Clustering: A Structure for Scalable Multiprocessor Operating System Design

RONALD C. UNRAU

*IBM Canada, 844 Don Mills Rd., North York, Ontario, M3C 1V7*

unrau@torolab2.vnet.ibm.com

ORRAN KRIEGER

*Department of Electrical and Computer Engineering, University of Toronto,  
Toronto, Canada M5S 1A4*

okrieg@eecg.toronto.edu

BENJAMIN GAMSA

*Department of Computer Science, University of Toronto, Toronto, Canada M5S 1A4*

ben@sys.toronto.edu

MICHAEL STUMM

*Department of Electrical and Computer Engineering, University of Toronto,  
Toronto, Canada M5S 1A4*

stumm@eecg.toronto.edu

*(Received March 1994; final version accepted June 1994.)*

**Abstract.** We introduce the concept of *hierarchical clustering* as a way to structure shared-memory multiprocessor operating systems for scalability. The concept is based on clustering and hierarchical system design. Hierarchical clustering leads to a modular system, composed of easy-to-design and efficient building blocks. The resulting structure is scalable because it 1) maximizes locality, which is key to good performance in NUMA (non-uniform memory access) systems and 2) provides for concurrency that increases linearly with the number of processors. At the same time, there is tight coupling within a cluster, so the system performs well for local interactions that are expected to constitute the common case. A clustered system can easily be adapted to different hardware configurations and architectures by changing the size of the clusters. We show how this structuring technique is applied to the design of a microkernel-based operating system called HURRICANE. This prototype system is the first complete and running implementation of its kind and demonstrates the feasibility of a hierarchically clustered system. We present performance results based on the prototype, demonstrating the characteristics and behavior of a clustered system. In particular, we show how clustering trades off the efficiencies of tight coupling for the advantages of replication, increased locality, and decreased lock contention.

**Keywords:** Operating systems, shared-memory multiprocessors, scalability, hierarchical design.

## 1. Introduction

Much attention has been directed towards designing “scalable” shared-memory multiprocessor hardware, capable of accommodating a large number of processors. These efforts have been successful to the extent that an increasing number of such systems exist [BBN 1989; Chaiken et al. 1991; Frank et al. 1993; Kuck et al. 1986; Lenoski et al. 1992; Oed 1993; Pfister et al. 1985]. However, scalable hardware can only be fully exploited and be cost-effective for general purpose use if there exists an operating system that is as scalable as the hardware.

An operating system targeting large-scale multiprocessors must consider both concurrency and locality. However, existing multiprocessor operating systems have been scaled

to accommodate many processors only in an ad hoc manner, by repeatedly identifying and then removing the most contended bottlenecks, thus addressing concurrency issues but not locality issues. Bottlenecks are removed either by splitting existing locks or by replacing existing data structures with more elaborate, but concurrent ones. The process can be long and tedious and results in systems that 1) have a large number of locks that need to be held for common operations, with correspondingly large overhead; 2) exhibit little locality; and 3) are not scalable in a generic sense, but only until the next bottleneck is encountered [Balan and Gollhardt 1992; Barach et al. 1990; Campbell et al. 1991; Chang and Rosenberg 1992; Peacock et al. 1992]. Porting an existing system designed for networked distributed systems is also unsatisfactory, because of the large memory requirements and consistency traffic caused by replicating code and data and because of the overhead and lack of locality in interprocessor interactions.

We believe that a more structured approach for designing scalable operating systems is essential and propose a new, simple structuring technique based on clustering and hierarchical system design called *hierarchical clustering*. This structuring technique leads to a modular system, composed from easy-to-design and efficient building blocks. It is applicable to a wide variety of large-scale shared-memory architectures and configurations. The structure maximizes locality, which is the key to good performance in *non-uniform memory access* (NUMA) systems, and provides for concurrency that increases linearly with the number of processors.

In this paper we describe hierarchical clustering and its advantages. We show how it can be applied in implementing a microkernel-based operating system and present the results of performance measurements obtained from a running prototype. We describe some of the lessons we learned from our implementation efforts and close with a discussion of our future work.

## 2. Target Environment

The design of an operating system is influenced both by the hardware platform it is responsible for managing and the demands of the applications it expects to serve. The internal structure of the operating system and the policies it enacts must accommodate these targets if it is expected to perform well. In this section we introduce the particular hardware and workload targets we assume for this paper and briefly describe how the operating system can accommodate them.

### 2.1. The Target Hardware

Our target hardware is the class of “scalable” shared-memory multiprocessors. Most current scalable shared-memory multiprocessors are based on segmented architectures, which allow the aggregate communication bandwidth to increase as new segments are added to the system (with little increase in hardware complexity). In these systems the memory is distributed across the segments, but remains globally accessible. Because the cost of accessing memory is a function of the distance between the accessing processor

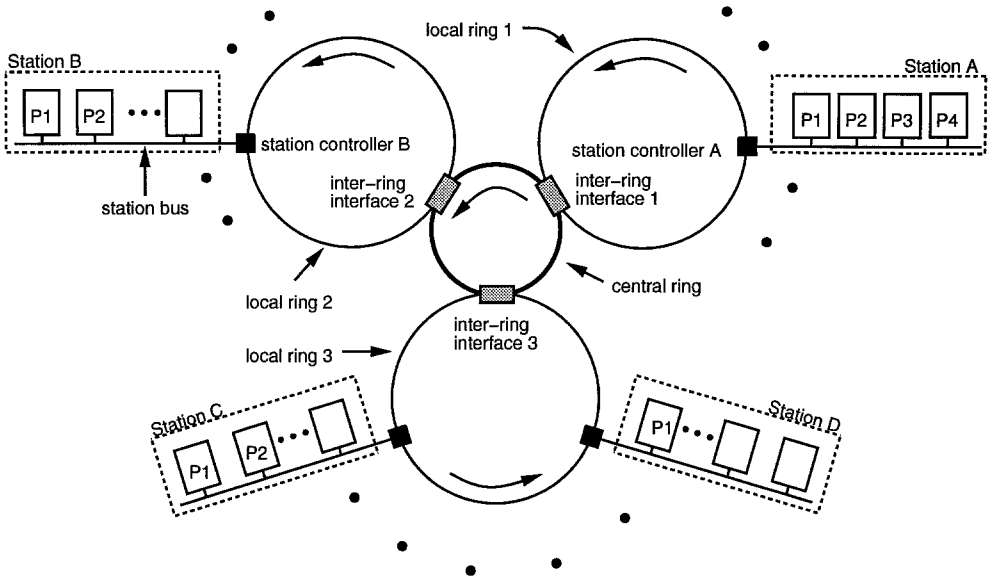


Figure 1. A high-level representation of the HECTOR multiprocessor. Each processor box in the system includes a CPU and cache, local memory, and I/O interface.

and the memory being accessed, they exhibit NUMA times. Examples of systems with segmented architectures include HECTOR [Vranesic et al. 1991], DASH [Lenoski et al. 1992], and the KSR1 [Frank et al. 1993]. The HECTOR system is depicted in Figure 1 and was used for the experiments discussed in Section 5.

Although the aggregate bandwidth of these scalable shared-memory machines increases with the size of the system, the bisectional bandwidth (from any one segment to the rest of the system) is typically much lower and limited. Both the higher cost of remote accesses and the limitations in bandwidth make it important to organize the location of data so that it is kept close to the accessing processors. In some cases, locality can also be improved by replicating data or moving it closer to the accessing processor, depending on how the data is being accessed. It should be noted that locality is important even for architectures based on cache-only memory (COMA) [Frank et al. 1993; Hagersten et al. 1992] and systems with hardware-based cache coherence [Chaiken et al. 1991; Lenoski et al. 1992], because operating system data is typically accessed with little temporal locality, yet shared at a fine granularity [Chen and Bershad 1993], leading to generally low cache hit rates.

## 2.2. *The Target Workload*

The workload of an operating system stems primarily from requests that come from application programs and will include, among other things, page-fault handling, I/O, interprocess communication, process management, and synchronization. The operating system is thus largely demand driven. The relative timing of the requests and the way they interact (by accessing common data structures) in the operating system depends on the type of applications currently running. The applications we expect to see running at any given time include a mix of sequential interactive programs, small-scale parallel applications, and very large parallel applications.

Each of these application classes places different demands on the operating system. The requests of an interactive, sequential application will be largely independent from the requests of the other applications, both with respect to the time they are issued and the operating system data they access. For example, the probability of two independent applications simultaneously accessing the same file is small. However, we can expect the load from these requests to grow linearly with the size of the system. To satisfy this load without creating a bottleneck, the number of operating system service centers (e.g., file servers or process dispatchers) must increase with the size of the system. Moreover, it is important to limit contention for locks that regulate access to shared data and to limit contention at physical resources, such as disks.

In contrast, the threads of a parallel application may well make simultaneous requests that cause the same system data to be accessed. This might happen, for example, after a program phase-change when each thread starts to access the same data file. When servicing these nonindependent requests, the accesses to the shared resources must be fair (or else starvation could result), and the latency of accessing these resources must degrade no worse than linearly with the number of concurrent requesters.

The operating system can help exploit the locality inherent in an application by scheduling communicating processes to neighboring processors, by directing I/O to nearby disks, and by placing application data close to where it will be accessed. Policies such as data replication or migration can also be used for both application memory and secondary storage management to increase access bandwidth, decrease latency, or both.

## 3. Hierarchical Clustering

Hierarchical clustering is a structuring technique based on clustering and hierarchical system design that we believe addresses most of the issues described in the previous section. Hierarchy is a well-known organizational framework, occurs frequently in nature [Simon 1985], and is often used in business, government, and the military. In these systems the purpose of the higher levels of the hierarchy is to manage or control resources in a coarse, strategic, and global fashion, while lower levels in the hierarchy control resources in a finer, more localized way. Communication occurs primarily between entities close to each other, and interaction between distant entities typically occurs through higher levels of the hierarchy. If the operations at the lowest levels are indeed

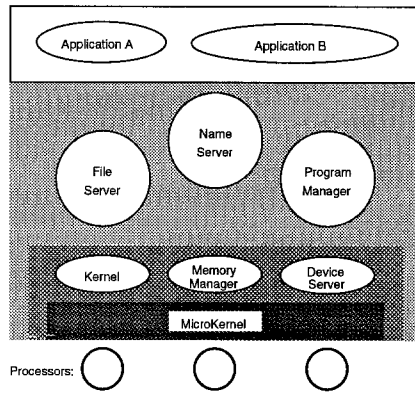


Figure 2. A single clustered system.

localized, then there can be much concurrency in the system. It is both the locality and concurrency aspects of these systems that allow them to scale to relatively large sizes.

A hierarchically clustered operating system is scalable for the same reason: It provides for locality and it provides for concurrency that increases with the number of processors. It avoids concurrency constraints imposed at the hardware level due to bus and network bandwidth limitations and imposed at the operating system level due to locks by 1) replicating data structures and by 2) exporting appropriate policies. The purpose of clusters is to provide for tight coupling and good performance for neighboring objects that interact frequently with one another. Hierarchical clustering also has some pragmatic advantages; for example, it simplifies lock tuning and is adaptable to different architectures and configurations.

In this section we first describe the basic structure and then argue why it is scalable. This is followed by a discussion of factors that affect the best choice of cluster size.

### 3.1. The Basic Structure

In a system based on hierarchical clustering the basic unit of structuring is a *cluster*. A cluster provides the functionality of an efficient, small-scale symmetric multiprocessing operating system targeting a small number of processors. Kernel data and control structures are shared by all processors within the cluster. In our implementation a cluster consists of a symmetric microkernel, memory and device management subsystems, and a number of user-level system servers such as a scheduler and file servers (Figure 2). Note, however, that clustering does not preclude the management of some resources on a per-processor basis within a cluster. For example, the performance of task dispatching can be improved if processes are maintained on per-processor rather than per-cluster ready queues.

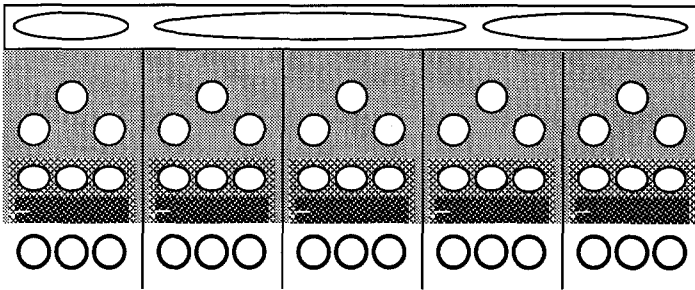


Figure 3. A multiple clustered system.

On a larger system, multiple clusters are instantiated such that each cluster manages a unique group of “neighboring” processing modules (including processors, memory, and disks). All major system services are replicated to each cluster so that independent requests can be handled locally (Figure 3). Clusters cooperate and communicate to give users and applications an integrated and consistent view of a single large system.

The basic idea behind using clusters for structuring is to use multiple computing components that are easy to design and hence efficient for forming a complete system. A single level of clusters can be expected to support moderately large systems (in the, say, 64- to 256-processor range). However, for larger systems, additional levels in the hierarchy will be necessary. In particular, while each cluster is structured to maximize locality, there is no locality in cross-cluster communication with only a single level of clusters. The logical next step is to group clusters into superclusters, and then to group superclusters into super-superclusters, in which each level in the clustering hierarchy involves looser coupling than the levels below it.

Hierarchical clustering incorporates structuring principles from both tightly coupled and distributed systems and attempts to exploit the advantages of both. By using the structuring principles of distributed systems, the services and data are replicated and situated so as to 1) distribute the demand, 2) avoid centralized bottlenecks, 3) increase concurrency, and 4) increase locality. On the other hand, there is tight coupling within a cluster, so the system is expected to perform well for the common case when interactions occur primarily between objects located in the same cluster.

Because of its hybrid structure, many design issues in a clustered system are similar to those encountered in distributed systems. For example, shared objects can be distributed and moved across clusters to increase locality, but then the search to locate these objects becomes more complex. Other objects may be replicated to increase locality and decrease contention, but then consistency becomes a critical issue.

There are also important differences between a clustered system and a distributed system that lead to completely different design tradeoffs. In a distributed system, for example, the hosts do not share physical memory, so the cost for communication between hosts is far greater than the cost of accessing local memory. In a (shared-memory) clustered system, it is possible to directly access memory physically located in other

clusters, and the costs for remote accesses are often not much higher than for local accesses. Additionally, demands on the system are different in the multiprocessor case, because of tighter coupling assumed by the applications. Finally, fault tolerance is a more important issue in distributed systems, in which the failure of a single node should not crash the entire system. In a shared-memory multiprocessor this type of fault tolerance is not (yet) a major issue (although clustering should provide a good foundation on which to address operating system design for fault containment when suitable hardware becomes available).

### 3.2. Scalability

Intuitively, scalability of an operating system refers to its capacity to provide increased service as the hardware resources (e.g., processors, memory, or disks) are increased. In queuing theoretic terms, an operating system can only scale if there are no resources within the system (e.g., queues and locks) that saturate as the system size is increased.

In structuring an operating system for scalability, it is important that application requests to independent physical resources be managed by independent operating system resources. For example, if two applications simultaneously request data from two different disks, concurrency in servicing the requests is maximized if the operating system structures that implement the service are completely independent.

More generally, a scalable operating system will possess the following properties [Unrau 1993]:

- *The parallelism afforded by the applications is preserved.*

The operating system workload can be expected to increase with the size of the system. Because most of the parallelism in an operating system comes from requests that are concurrently issued by multiple application threads, the number of service centers must increase with the size of the system, as must the concurrency available in accessing data structures.

With hierarchical clustering, each cluster provides its own service centers and has its own local data structures. Application requests are always directed to the local cluster. Hence, the number of service points to which an application's requests are directed grows proportional to the number of clusters spanned by the application. As long as the requests are to independent physical resources, they can be serviced in parallel.

- *The overhead for independent requests is bounded by a constant [Barak and Koratzsky 1987].*

If the overhead of a service call increases with the number of processors, then the system will ultimately saturate. For this reason the demand on any single operating system resource cannot increase with the number of processors. With hierarchical clustering, there are no systemwide ordered queues and, more generally, no queues whose length increases with the size of the system. Because the number of processors in a cluster is fixed, fixed-sized data structures can be used within a cluster to bound the overhead of accessing these structures.

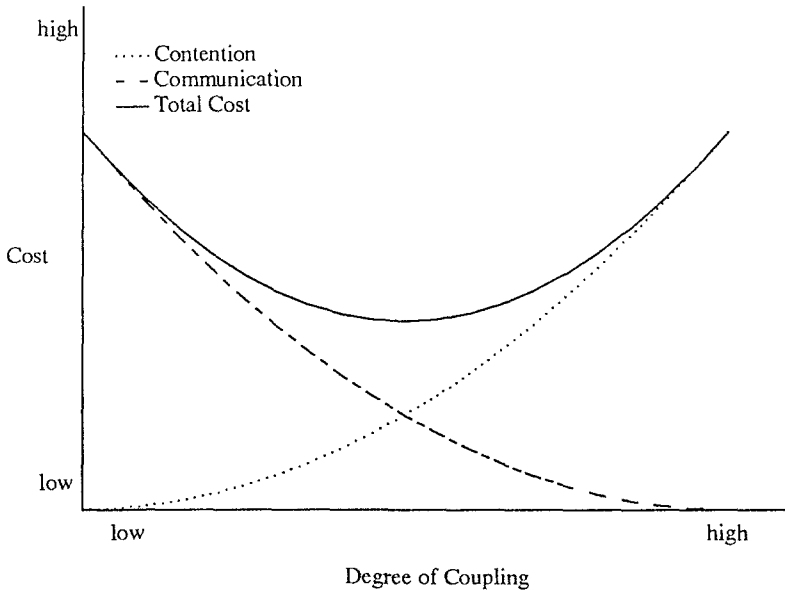


Figure 4. Cost versus coupling in a large-scale system.

- *The locality of the applications is preserved.*

The locality of an application is inherent in the algorithm; the operating system cannot create locality where none existed before, but it must preserve and exploit the locality that is present if the application is to perform well. Hierarchical clustering exploits locality by co-locating objects that need to interact frequently in the same cluster (or supercluster if the interactions cross cluster boundaries) and by always first directing requests from the application to the local cluster. For example, descriptors of processes that interact frequently are forced to be close to each other, and memory mapping information is located close to the processors that need it for page-fault handling. Hierarchical clustering also provides a framework for enacting policies that increase locality in the applications' memory accesses and system requests. For example, policies can attempt to run the processes of a single application on processors close to each other, place memory pages near the processes accessing them, and attempt to direct file I/O to nearby devices.

### 3.3. Cluster Size Tradeoffs

The above requirements for scalability can be largely met by using a distributed systems approach, in which system services and structures are replicated onto each individual processing module; that is, by using a cluster size of one and a single level of hierarchy. The many large-scale distributed (i.e., networked) systems in existence today clearly demonstrate this. However, using a cluster size of one on a shared-memory multiprocessor will, for several reasons, result in greater overhead than using larger cluster sizes. First, the smaller the cluster size, the more intercluster operations are incurred, and intercluster



operations are more costly than intracluster operations. Second, there is less locality in intercluster communication, which further increases its cost. Third, it is necessary to keep the state on each cluster consistent with that on the other clusters. Finally, a cluster size of one results in greater memory requirements due to the replication of much of the operating system's code and data.

A cluster size larger than one reduces the consistency and space overhead and reduces the number of intercluster operations required. However, the size of a cluster cannot become arbitrarily large without reorganizing the internal cluster structure, because the probability of contention for shared data structures will increase with the number of processors in the cluster, ultimately resulting in saturation.

These tradeoffs in the choice of cluster size are depicted abstractly by the graph in Figure 4. The graph plots cost (including overhead) against the degree of coupling or cluster size for some fixed number of processors. The dotted curve in the figure represents access overhead and contention, which increases as the cluster size becomes larger because the probability of active resource sharing increases, while locality within the cluster decreases. The dashed curve represents the portion of overhead due to remote, intercluster communication; in general, fewer intercluster operations are necessary when the size of the cluster is larger. The solid curve is the sum of the contention and communication costs and suggests that there is some degree of coupling for which the overall cost is minimized. In Section 5 we show that this cost behavior can actually be observed in real systems.

One of the benefits of hierarchical clustering is that system performance can be tuned to particular architectures, configurations, and workloads by adjusting the size of the clusters. The architectural issues that affect the appropriate cluster size include the machine configuration, the local-remote memory access ratio, the hardware cache size and coherence support, and the topology of the interconnection backplane. On hierarchical systems such as Cedar [Kuck et al. 1986], DASH [Lenoski et al. 1992], or HECTOR [Vranesic et al. 1991], the lowest-level operating system cluster might correspond to a hardware cluster. On a local-remote memory architecture, such as the BBN Butterfly [BBN 1988], a smaller cluster size (perhaps even a cluster per processor) may be more appropriate; in this case, clustering can be viewed as an extension of the fully distributed structuring sometimes used on these machines.

Although architectural characteristics are important, they are not the only factor that determines the best cluster size for a particular system. For example, some systems have relatively large hardware cluster sizes, like the KSR1 [Frank et al. 1993] with 32 processors per ring. If one were to choose a cluster size of 32 for this machine, it would be difficult to prevent the locks in the operating system from saturating.

Choosing the correct locking granularity for a parallel system is generally difficult. A large granularity restricts concurrency, but a fine granularity increases the length of the critical path. Hierarchical clustering allows the locking granularity within a cluster to be optimized for the size and concurrency requirements of the cluster [Unrau et al. 1994]. For example, the message-passing code in our implementation uses only a small number of large-grained locks, adding little overhead to the critical path but providing sufficient concurrency for the expected demand given our commonly used cluster size. Fine-grained

locking is used when a high degree of concurrency is required or if only a small number of locks need to be acquired in the critical path.

#### 4. Implementation

We have implemented a microkernel-based operating system called HURRICANE, structured according to the principles of hierarchical clustering. The prototype does not as yet support the superclusters we envisage are necessary for large systems. Instead, we have focused initial efforts on structuring within and across first-level clusters. At present, the system is fully functional, and in this section we show how clustering has affected the structure of its primary components.

In order to make a clustered operating system successful, it is important to keep the following issues in mind. Applications should be shielded from the fact that the operating system is partitioned into clusters; to users and applications the system should, by default, appear as a single, large, integrated system. Moreover, the complexity introduced by separating the system into clusters must be kept reasonably low, since clusters were introduced as a structuring mechanism primarily to reduce overall complexity, not increase it. Finally, intercluster communication must be efficient; while we expect most interactions to be local to a cluster, those interactions that must access remote clusters should not be overly penalized.

The main design issue in implementing a clustered system (as opposed to a traditionally structured system) is the choice in strategy for managing the data and the choice in mechanism for communication. For each type of data, one must decide whether data of that type should be replicated or whether there should be only one copy. If there is only one copy, there is a choice of placing it statically or allowing it to dynamically migrate to where it is being accessed. If data is replicated, then one must decide how to keep it consistent. The above choices will be influenced primarily by the expected access behavior. For example, one would not want to replicate actively write-shared data.

Intercluster communication is needed 1) to access remote data that is not replicated, 2) to replicate and move data, and 3) to maintain data consistency. Three different mechanisms can be used for intercluster communication [Chaves et al. 1993]. First, shared-memory access is always a viable option on a shared-memory multiprocessor, particularly for lightweight operations that make only a few references. For instance, our implementation uses shared memory for table lookup to locate objects such as page descriptors. However, we prefer to minimize cross-cluster shared-memory accesses because of the reduced locality in these accesses (and hence the increased probability of contention if used indiscriminately). There are also other problems with accessing remote resources using shared memory. With spin locks, secondary effects such as memory and interconnect contention can result in a greater than linear increase in cost as the number of contending processors grows [Anderson 1990; Mellor-Crummey and Scott 1991]. When only processors belonging to the same cluster contend for a spin lock, these secondary effects can be bounded.

A second possibility is to use *remote procedure calls* (RPCs). In our system RPCs are implemented using remote interrupts and are used for kernel-kernel communication. The

interrupted cluster obtains the call information and arguments through shared memory directly, after which data accesses required to service the request can be local to the interrupted cluster. The savings gained from the local data accesses must amortize the cost of the interrupt (i.e., the cycles that are stolen and the register saving and restoring).

Finally, HURRICANE supports message passing that allows user-level processes to send messages to other processes regardless of their location in the system. This provides a simpler form of cross-cluster communication for higher-level servers that are replicated across the clusters of the system.

#### 4.1. *The HURRICANE Kernel*

The HURRICANE kernel is responsible for process management and communication and is similar in structure to the V kernel [Cheriton 1988]. Most operations on processes involve the queuing and dequeuing of process descriptors. For example, when a process sends a message to another process, its descriptor (containing the message) is added to the *message queue* of the target descriptor. Other queues include a *ready queue* for processes that are ready to run and a *delay queue* for processes waiting for a time-out. Each HURRICANE cluster has its own instance of kernel data structures. Hence the number of queues increases with the size of the system, so the average demand on each queue should remain constant as the system grows. Moreover, the process descriptor table and kernel queues are local to each cluster, preserving the locality for the common case when the request can be handled entirely within the cluster. A process descriptor moves with the process if it migrates across cluster boundaries so that operations on the process descriptor are always local to the cluster on which the process resides and so that queue traversals are also always local to the cluster.

A home encoding scheme is used to allow a process to be located anywhere in the system with a constant number of operations, independent of the number of processors. When a process is created, it is assigned a process identifier within which the ID of the *home* cluster is encoded. Usually, a process remains within its home cluster, but if it migrates, then a record of its new location (a *ghost* descriptor) is maintained at the home cluster. Operations that require access to the newly located process look first in the home cluster, find the ghost, and use the locator information in the ghost to determine the current cluster in which the process resides.

Intercluster message passing is implemented using remote procedure calls that interrupt the target cluster to 1) allocate a local message retainer (similar to a process descriptor) in the target cluster, 2) copy the message and other information into the retainer by directly accessing the process descriptor in the source cluster, and 3) add the message retainer to the message queue of the target process descriptor. Note that by using local message retainers in the target cluster (as opposed to a global pointer to the process descriptor in the source cluster), the message queue will always be an entirely local data structure.

In the case of the HURRICANE kernel, having the cluster size smaller than the size of the system has three key advantages:

1. It localizes the kernel data structures that need to be accessed;

2. it reduces the number of process descriptors that must be managed within a cluster, thus reducing the average length of clusterwide queues; and
3. systemwide, it increases the number of queues and tables (since each cluster has local instances), reducing demand on the resources and increasing concurrency.

On the other hand, having a cluster span of more than one processor reduces the amount of intercluster communication and makes it easier to balance the load of the processors.

#### 4.2. *Protected Procedure Calls*

With a microkernel multiprocessor operating system, the majority of IPC (interprocess communication) operations are required solely to cross address space boundaries (i.e., for client-server interactions) and not to cross processor boundaries. The HURRICANE message-passing facility described above is optimized for interprocessor communication, whereas another facility, called *protected procedure calls* (PPCs), is optimized for interaddress space communication.

The PPC facility uses the same structuring principles as hierarchical clustering but applied on a per-processor basis. As a result, it accesses only local data, requires no access to shared data, and requires no locking, so end-to-end performance is comparable to the fastest uniprocessor IPC times. Yet because of the concurrency provided, this performance can be sustained independent of the number of concurrent PPC operations, even if all such operations are directed to the same server. Moreover, the PPC facility ensures that local resources are made available to the target server so that any locality in client-server interactions can be exploited by the server [Gamsa et al. 1993].

In our implementation of PPCs, separate worker processes in the server are created as needed to service client calls on the same processor as the requesting client. On each call they are (re)initialized to the server's call handling code, which has the same effect as an up-call into the service routine. This results in a number of benefits: 1) Client requests are always handled on their local processor, 2) clients and servers share the processor in a manner similar to hand-off scheduling, and 3) there are as many threads of control in the server as client requests.

The key benefits of our approach all result from the fact that resources needed to handle a PPC are accessed exclusively by the local processor. By using only local resources during a call, remote memory accesses are eliminated. More importantly, since there is no sharing of data, cache coherence traffic can be eliminated. Since the resources are exclusively owned and accessed by the local processor, no locking is required (apart from disabling interrupts, which is a natural part of system traps).

Although call handling is managed at the local processor level, the creation and destruction of communication endpoints are managed at the cluster level (in order to synchronize these events across the cluster), requiring cross-cluster communication only for those servers that span multiple clusters. Such operations are sufficiently rare in our system that the performance impact is negligible.

### 4.3. Memory Management

Each HURRICANE virtual address space defines the application-accessible memory resources as a set of non-overlapping regions, each of which describes the attributes of a contiguous sequence of virtual memory pages. Each region is directly associated with a corresponding file region so that accesses to the memory region behave as though they were accesses to the file region directly. This relationship is known as a single-level store, or mapped-file, memory abstraction [Unrau 1993]. Within a cluster, address spaces are represented by an address space descriptor and an associated balanced binary tree of region descriptors. The region descriptors maintain the virtual limits and mapping attributes of each address space region.

For programs that span multiple clusters, the address space and region descriptors are replicated to the clusters as needed. To maintain the integrity of these replicated structures, all modifications are directed through the *home cluster*, which serves as a point of synchronization. Currently, the home cluster is designated as the cluster on which the program was created. Since the address space and region descriptors are accessed for reading far more often than they are modified, replication of these structures improves locality of access in the common case.

In large systems, applications that span many clusters could cause the memory manager on the home cluster to become contended, because it must both replicate the data structures and distribute modifications to them. Additional levels in the hierarchy could be used to improve both cases. On first access, a source replica for the data structure could be obtained from the local supercluster or super-supercluster rather than the home cluster. A spanning tree could then be used to distribute modifications to all the replicas.

Each physical page in a cluster has a page descriptor associated with it. This structure identifies the physical address of the page, its state, and the file block stored there. All pages that contain valid file data are kept in a *page cache*, which is a hash table indexed by file block and is used to determine whether a file block is already resident in the cluster. Unreferenced pages are placed on a free list, but remain in the page cache in case the file block is referenced again in the near future. These structures are always local to a cluster and need not be replicated or moved, since they represent local physical resources. To reduce the amount of cross-cluster communication, representative page descriptors are used to identify pages that are accessed remotely. This keeps all the data structures local to the cluster and allows the results of a remote page search to be subsequently used by other processors within the cluster.

The HURRICANE memory manager supports page replication and migration across clusters. This can reduce access latency and contention for some applications [Bolosky et al. 1989; Cox and Fowler 1989; LaRowe et al. 1991]; however, the overhead of these policies must be amortized to realize a net gain in performance. Hierarchical clustering provides a framework for enacting paging policies, in which pages are shared within a cluster to reduce overhead, but are replicated and moved across clusters to reduce access latency and memory contention.

A simple directory mechanism is used to maintain page consistency between clusters. There is one directory entry for each valid file block currently resident in memory; it

identifies which clusters have copies of the page. The directory is distributed across the processors of the system (thus spanning all clusters), allowing concurrent searches and balanced access demand across the system. Pages that are actively write-shared are not replicated, but instead are accessed remotely.

#### 4.4. The File System

File I/O in HURRICANE is provided largely outside of the kernel; only the device drivers are in the kernel. The HURRICANE File System (HFS) implementation makes use of three user-level system servers [Krieger and Stumm 1993; Krieger 1994]. The *Name Server* manages the HURRICANE name space. The *Open File Server* (OFS) maintains the file system state for each open file and is largely responsible for authenticating application requests to lower-level servers. The *Block File Server* (BFS) controls the disks, determines the target for each I/O request, and directs the request to a corresponding device driver.

The *Alloc Stream Facility* (ASF) is a file system library in the application address space that translates application file `read` and `write` requests into accesses to a mapped file region [Krieger 1994; Krieger et al. 1994]. On the first access to a page in the region, the process will page-fault. If the data is not available in the file cache, then the fault is translated by the memory manager into a `read_page` request to the BFS.

There is a separate instance of all file system servers on each cluster. Client/server interactions within a cluster are directed to servers local to the client. For example, the memory manager directs all page-fault requests to its local BFS. Only requests between peer HFS servers (for example, BFSs in different clusters) may cross cluster boundaries. Cross-cluster communication between file system servers occurs through explicit message passing, not through shared memory.

When an application spans multiple clusters, or when programs on different clusters access the same files, then the file system state for shared files is replicated to the cluster servers as they are needed. When the state is first required by a server, it locates some other server that has the state to obtain a local copy of the state. Currently, this is done by directing a request to the server in the cluster identified by the persistent token of the file. To maintain the integrity of these replicated structures, all modifications are directed through the *home cluster*, which serves as a point of synchronization. In our current design the home cluster is the cluster where the file was first opened. Depending on the particular type of state being accessed, modifications result in either update or invalidate messages being sent to the servers on the clusters that have a copy.

In addition to providing a structured way to achieve scalability for the file system servers, clustering provides a framework for invoking policies that localize application I/O requests and balance the load across the disks of the system. Files are created on disks that are in clusters near the requesting applications. If a file is replicated to multiple disks, the file system can use hierarchical clustering to determine the nearest replica of the file and direct read requests to that replica. Similarly, for some kinds of distributed files the file system can distribute write requests to disks in clusters near the page frame that contains the dirty page. In the case of load balancing, HFS maintains both long- and

short-term statistics on the aggregate disk load of the clusters and superclusters. This information can be used in order to direct I/O requests to the least loaded clusters. This information can also be used when a new file is created so that the file is created on disks in lightly loaded clusters.

#### *4.5. Scheduling*

One goal of the scheduling subsystem is to keep the load on the processors (and possibly other resources, such as memory) well balanced. The scheduling decisions in HURRICANE are divided into two levels. Within a cluster the load on the processors can be balanced at a fine granularity through the dispatcher (in the microkernel). This fixed scope limits the load on the dispatcher itself and allows local placement decisions on different clusters to be made concurrently. Cross-cluster scheduling is handled by higher-level scheduling servers, which balance the load by assigning newly created processes to specific clusters and by moving existing processes to other clusters. The coarser granularity of these movements permits a low rate of intercluster communication between schedulers.

A second goal in a hierarchically clustered system is to enact process placement policies that exploit the locality of application programs. For parallel programs with a small number of processes, all of the processes should run in the same cluster. For larger-scale parallel programs that span multiple clusters, the number of clusters spanned should be minimized. These policies are motivated by studies that have shown that clustering can noticeably improve overall application performance [Ahmad and Ghafoor 1991; Brecht 1993; Feitelson and Rudolph 1990; Zhou and Brecht 1991].

### **5. Experimental Results**

An important aspect of our implementation is that the cluster size can be specified at boot time. This allows us to explore the tradeoffs between contention and communication overhead and to quantitatively evaluate our locality management policies. This section reports the results of several experiments designed to evaluate hierarchical clustering for both synthetic stress tests and real applications. All experiments were run on a fully configured version of HURRICANE with all servers running, but with no other applications running at the time.

The experiments were run on a 16-processor prototype of a hierarchical NUMA architecture, HECTOR (see Figure 1). Although a 16-processor system is relatively small, the impact of contention and the NUMA times of the system are clearly evident and will only get worse in larger systems with faster processors. Hence we believe that the general trends exhibited by the results of our experiments can be extrapolated to larger systems. Also, our results show the basic costs for a clustered system of inter- and intracluster operations, which are largely independent of the system size.

The prototype system used for experimentation uses Motorola 88100 processors running at 16.67 MHz, with 16-Kbyte data and instruction caches (88200), and a local portion

of the globally addressable memory. At the lowest level, processors are connected by a bus. These lowest-level units, called *stations*, are then connected together by a high-speed ring, with multiple rings connected together by higher-level rings, and so on. The particular hardware configuration used in our experiments consists of four processing modules per station and four stations connected by a ring. By comparison with other NUMA systems [BBN 1989; Chaiken et al. 1991; Frank et al. 1993; Kuck et al. 1986; Lenoski et al. 1992; Pfister et al. 1985] the cost of accessing remote memory on HECTOR grows relatively slowly as the system is scaled. The time to load a cache line grows from 19 cycles for local memory, to 29 cycles for a different memory module in the same station, up to 34 cycles when the memory module is in a different station. The HECTOR prototype on which we ran these experiments does not support hardware cache coherence, but page-based cache coherence is supported by the HURRICANE memory manager.

### 5.1. Basic Operations

We now present some of the basic costs of the system and compare intracluster versus intercluster times for some common operations.

Protected procedure calls (PPCs) are only supported for local (on-processor) communication, and hence their cost is independent of clustering. A null-PPC from a user application to a user-level server is 32.4  $\mu\text{s}$  and drops to 19.2  $\mu\text{s}$  when directed to a kernel-level server (both are round-trip times). The performance of PPCs remains constant irrespective of the number of processors making concurrent requests.

*Remote procedure calls* are used for cross-cluster communication within the operating system kernel. The cost of a cross-cluster null-RPC in our implementation is 27  $\mu\text{s}$ , which includes the time to interrupt the target processor and the time to save and subsequently restore its registers.

The cost to handle a read page-fault for an in-core page increases from 128  $\mu\text{s}$  within a cluster to 243  $\mu\text{s}$  across clusters. The additional overhead in the latter case is primarily due to the overhead of replicating the page descriptor and the extra directory lookup needed to locate the remote copy of the page. Subsequent faults by other processors on the remote cluster incur only the in-cluster cost.

The cost of a 32-byte send-response message transaction between two user-level processes increases from 328  $\mu\text{s}$  within a cluster to 403  $\mu\text{s}$  across clusters. Here the extra overhead is due to two RPCs needed to communicate between clusters and the allocation, initialization, and subsequent deallocation of message retainers at the target cluster.

### 5.2. Synthetic Stress Tests

To study the effects of clustering on throughput and response time, we designed several synthetic tests to stress the concurrency limits of the system. A message-passing test is used because message passing is a key means of communication between the servers of a microkernel operating system. A test of the performance of page faults to in-core



*Table 1.* The cluster configurations used for the synthetic stress tests and application benchmarks. **CSize-16** corresponds to a tightly coupled system; **CSize-1** corresponds to a fully distributed system.

<b>CSize-1</b>	-	16 clusters with 1 processor each
<b>CSize-2</b>	-	8 clusters with 2 processors each
<b>CSize-4</b>	-	4 clusters with 4 processors each
<b>CSize-8</b>	-	2 clusters with 8 processors each
<b>CSize-16</b>	-	1 cluster with 16 processors

pages (soft page faults) is used, since such soft page faults are fairly common, as they are required to support page-level cache coherence and page migration and replication, both important for performance in a NUMA multiprocessor. There are two groups of tests, both containing the message-passing and page-fault test. The first group exercises the system for simultaneous operations to *independent* resources; the second group shows the performance for concurrent operations to *shared* resources.

The tests were run on each of the cluster configurations of Table 1. Note that the hardware configuration does not change and is always four processors to a station and that the clusters are assigned to the hardware stations in a natural way.

The cluster configurations correspond to different degrees of coupling and therefore allow us to explore the communication/contention tradeoffs alluded to in Figure 4. We expect **CSize-1**, in which each processor has its own set of operating system data structures, to perform well for coarse-grained parallel or independent applications, but to suffer from high communication costs for tightly coupled applications. At the other end of the spectrum, **CSize-16** has one set of kernel data structures that are shared by all processors in the system. For this configuration we expect contention for the shared data structures to be a significant portion of the response times.

### 5.2.1. Independent Operations

Figure 5 is a schematic depiction of the programs that stress simultaneous operations on independent resources. In the message-passing test, neighboring pairs of processes repeatedly exchange 32-byte messages. In the page-fault test,  $p$  processes repeatedly fault on a separate region of local memory. The pages of the region are unmapped between separate iterations of the loop. Figure 6 shows the response times and normalized throughput of the tests of Figure 5. Normalized throughput is calculated as  $X(p) = T(1)/T(p) \cdot W(p)/W(1)$ , where  $T(p)$  is the response time for  $p$  processes, and  $W(p)/W(1)$  reflects the amount of work done relative to the single operation case. The smallest cluster size shown for the message-passing test is **Csize-2**, since it is the first configuration for which neighboring processes reside in the same cluster. In these

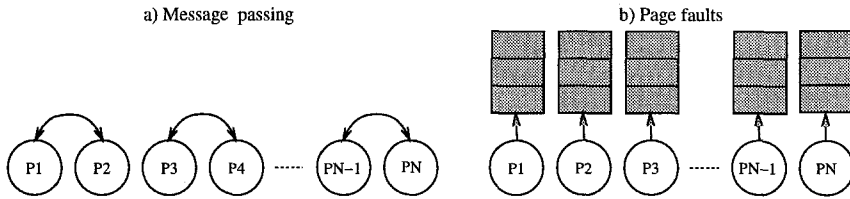


Figure 5. A schematic depiction of the programs that stress simultaneous operations on independent resources. In the message-passing test (a), neighboring pairs of processes repeatedly exchange 32-byte messages. In the page-fault test (b),  $p$  processes repeatedly fault on a separate region of local memory.

tests the  $p$  processes are assigned to separate processors in a way that minimizes the number of clusters spanned. That is, the processes are first added to one cluster until each processor in the cluster has a process, before adding processes to the next cluster. This corresponds to a scheduling policy that attempts to assign related processes to the same cluster whenever reasonable.

Figure 6 (a) and (b) shows that for both tests, the response time increases until  $p$  is equal to the number of processors in a cluster, at which point lock contention within the cluster is maximized. For example, in the page-fault experiments the address space structure and the top levels of the balanced binary tree of regions become the contended resources. As the processes “spill over” to the next cluster, the response times drop because the new processes experience no contention and thus lower the overall average response times.

Observe that the response times for single operations (the lower left corner of the figure) increase slightly as the cluster size increases. This is due to NUMA effects in the larger clusters. Since there is only one operation in progress, there is no lock contention, but the probability that a particular data structure is remote increases with cluster size, since the shared kernel resources are distributed across all the processors of the cluster. Thus the latency of kernel data accesses increases with cluster size.

Not surprisingly, the tests reveal that the smallest cluster size for each test has the best performance. Locality is maximized in this case, and since the operations are independent, no cross-cluster operations are required. For large cluster sizes the locks for common data structures suffer increasing contention, and hence response times increase correspondingly.

In summary, the results show that when clusters are large, contention can severely affect performance, and even when there is not much concurrency, the loss of locality will adversely affect performance to some degree. While these results were obtained from a small system, it is clear that both effects would be stronger in a larger system. The results would also be stronger on systems with more pronounced NUMA characteristics than HECTOR has (e.g., [BBN 1989; Chaiken et al. 1991; Frank et al. 1993; Lenoski et al. 1992]).

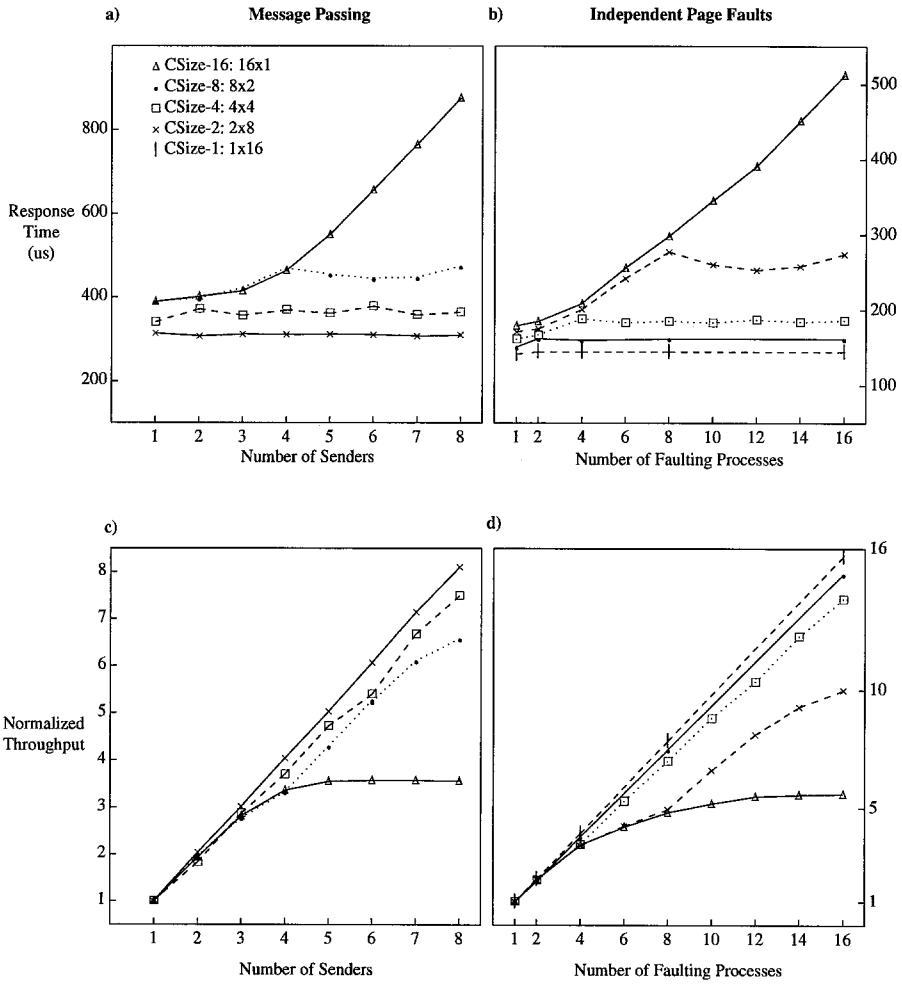


Figure 6. The response times and normalized throughput for simultaneous independent operations: (a) and (c) message passing; (b) and (d) page faults.

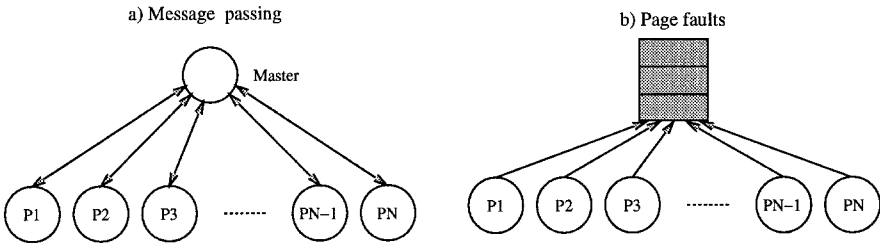


Figure 7. A schematic depiction of the programs that stress simultaneous access to shared resources. In the message-passing test (a),  $p$  children repeatedly send 32-byte messages to a single parent. In the page-fault test (b),  $p$  children simultaneously write to the same small set of shared pages.

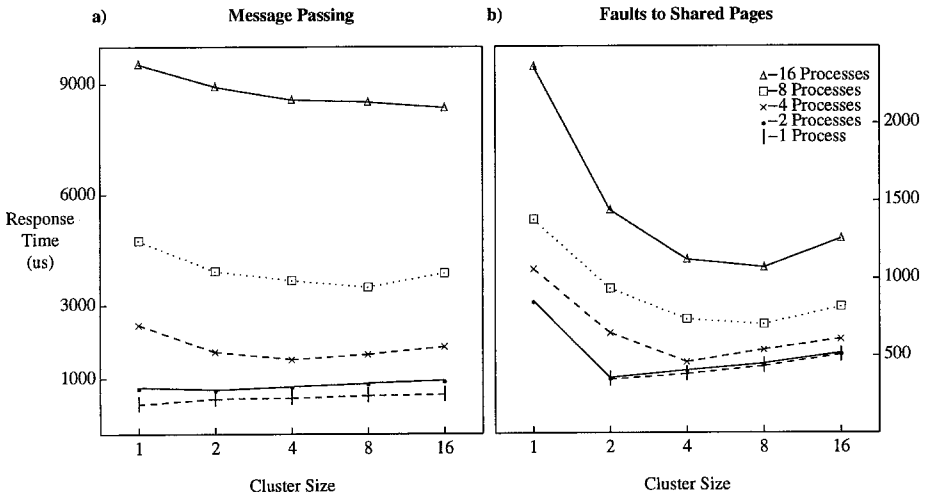


Figure 8. The response times for simultaneous operations to shared resources: (a) message passing; (b) page faults.

### 5.2.2. Operations to Shared Resources

Figure 7 is a schematic depiction of the programs that stress simultaneous access to shared resources. In the message-passing test,  $p$  children repeatedly send 32-byte messages to a single parent. In the page-fault test a master maps a small set of pages for reading, and then  $p$  children simultaneously write to the shared pages. The pages are unmapped between successive iterations of the test. Figure 8 shows the response times for the tests of Figure 7, plotted as a function of cluster size. The different curves represent different numbers of active children.

Figure 8 shows that an intermediate degree of coupling has the best performance for both tests. Also, the minimum response time for  $p$  processes occurs when the cluster

contains  $p$  processors, which suggests that it may be beneficial to match the cluster size to the application communication patterns. When the cluster size is smaller, the cost of cross-cluster operations is incurred; if the cluster size is larger, locality suffers.

The exception to the matching observation occurs for the 16-process page-fault test, for which **CSize-8** achieves the best performance. This is due to contention for shared resources, which dominates for the large cluster sizes. For small cluster sizes, the cost of the intercluster communication is more dominant.

Depending on the expected workload, there are two ways these results might be extrapolated to systems larger than the prototype. Consider first the case of a workload consisting primarily of large numbers of small-scale parallel applications (16 processes or under) running on a significantly larger system (several hundred processors). The results in this section show that as the cluster size is increased beyond the size of the application, performance degrades due to a loss of locality in the operating system. This trend is likely to continue in the future, when we expect even higher remote memory access costs, suggesting that clusters should remain relatively small (but larger than the average parallelism of the applications) in order to efficiently support this type of workload.

The second case to consider is a workload of large parallel applications. Here the results are not as clear cut. Extrapolating the message-passing test would lead to the conclusion that the cluster size should grow with the applications. However, the page-fault test demonstrates the danger of increasing the cluster size too much, even for tightly coupled applications. In particular, for the 16-process case, the optimal cluster size is 8, with a sharp rise in cost for cluster size 16. On the other hand, a cluster size significantly smaller than the application also results in poor performance. We expect from these results that increasing the cluster size will result in a diminishing improvement in performance until ultimately an increase in cluster size will decrease performance, suggesting an intermediate cluster size should be chosen even for tightly coupled large-scale applications.

The page-fault experiment of Figure 7 can be used to illustrate how additional levels of hierarchy can help performance on larger systems. To support coherence of read-shared pages at the time of a write, processors that are accessing the page must be notified. This is done using a spanning tree, in which a single processor in each cluster is notified and that processor is then responsible for notifying all other processors in its cluster. On larger systems with superclusters, this spanning tree can be extended to further improve concurrency.

### 5.3. Applications

The synthetic stress tests of the previous section were designed to exercise a single aspect of the system, namely simultaneous independent operations and simultaneous operations to shared resources. However, real applications seldom incur only one or the other and, in fact, exhibit both types.

We now examine the performance of three parallel applications as a function of cluster size. The applications are SOR, a partial differential equation solver; Matrix Multiply; and 2D-FFT, which calculates the Fourier transform of a two-dimensional array of data.

All three programs are of the data-parallel, or SPMD, class of applications. However, the data access patterns of the applications are very different so that each stresses a different aspect of the system.

Figure 9 shows the page-fault behavior of the three applications when four worker processes are used. In the figure the Y axis indicates the virtual page number (virtual address divided by the page size), and the X axis gives the time in milliseconds. Each "point" represents a page-fault event for the corresponding process and is plotted at the virtual address and time for which the fault occurred. Page-fault events are caused either by the process mapping the page into its address space on its processor or by writes trapped in order to support page-level coherence. For each of the applications the computation is partitioned along equal strips of the result matrix. The page-fault behavior is inherent in the algorithm and does not change (except for time-dependent accesses) as the operating system configuration is varied. However, the cost of handling the page faults changes with the configuration, depending on whether the fault is to a shared page or not.

The SOR application exhibits concurrent faults to primarily independent pages, since the only shared data is along the shared edges of the strips of the matrix. Figure 9 shows that there is little paging activity after the initial iteration of the algorithm, so we expect this application to perform well for small cluster sizes.

In contrast, the Matrix Multiply and 2D-FFT programs have phases of intense concurrent accesses to both independent and shared pages. For Matrix Multiply ( $C = A \times B$ ), the program starts with a period of concurrent accesses to independent pages, as the workers initialize the matrices in parallel. The computation starts at time 200 ms and proceeds along the rows of  $C$  (and  $A$ ). The accesses to these matrices are seen as the slowly increasing diagonals of fault instances in Figure 9. To compute a single row of  $C$ , each worker must read all the columns of  $B$ , which is seen as the sharp diagonals of faults in the figure. This access behavior stresses the memory manager's ability to handle simultaneous faults to shared pages.

The 2D-FFT application also begins with a flurry of faults to independent pages as the matrix and a reference matrix to check the answer are initialized by the workers in parallel. The forward transform begins along the rows at time 200 ms. At time 400 ms there is a period of intense coherence faults as the rows that had been accessed by a single worker are now accessed by all the workers to perform the column transform. The row-column fault pattern is repeated for the reverse transform, which begins at time 700 ms.

Figure 10 shows the response times for the three applications (on larger data sets than those shown in Figure 9) as the cluster size is varied from 1 to 16. The response times are plotted normalized to the response time of **CSize-4**. We are interested primarily in the general trends represented by each application, and hence the results for each application should be interpreted qualitatively and independently from each other. All the programs were run with 16 processes on 16 processors; only the cluster configurations were varied. The response times reported are for the entire program, which includes the time to create, schedule, and destroy the worker processes as well as the time to initialize any data structures. As a result, the times reported include the overhead due

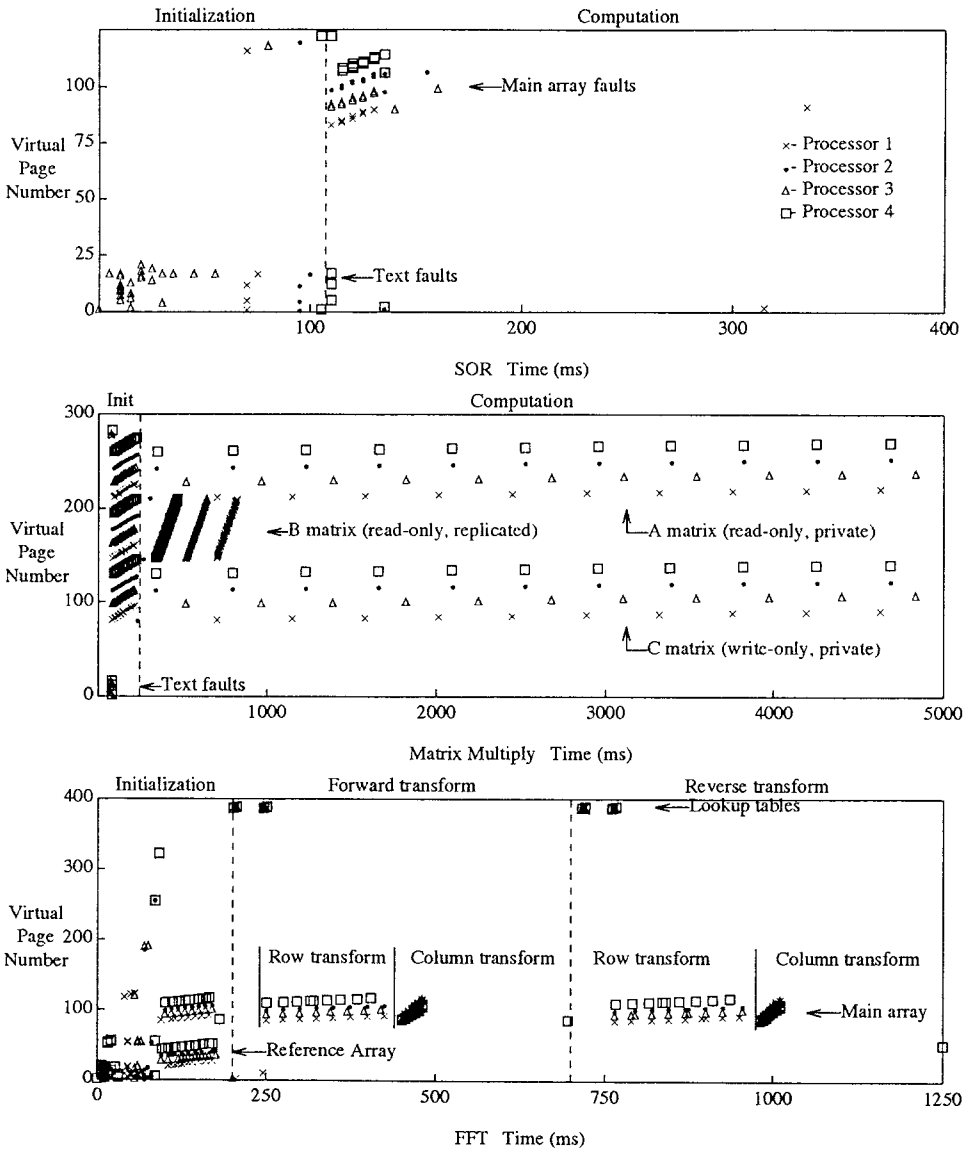


Figure 9. The page access and coherence faults of the three application programs (using four workers). The Y axis indicates the virtual page number (virtual address divided by the page size). The X axis gives the time in milliseconds. Each "point" represents a page-fault event for the corresponding process. Page-fault events are caused either by the process mapping the page into its address space on its processor or by writes trapped in order to support page-level coherence.

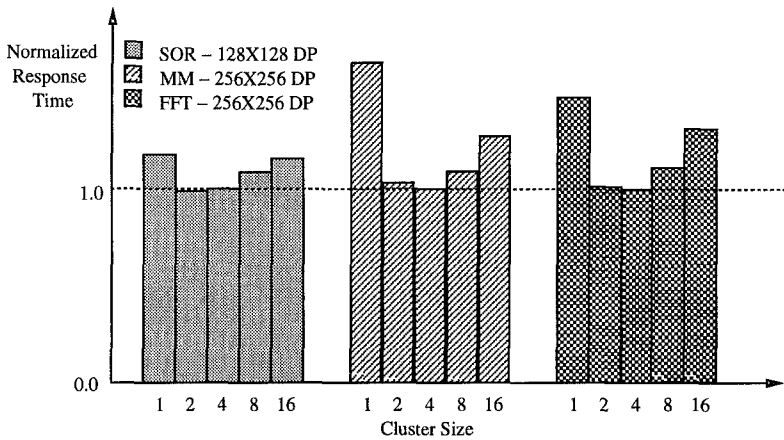


Figure 10. The normalized response times of the application suite for different data set sizes and cluster configurations.

to both process and memory management. The algorithms (and problem sizes) used for these applications were chosen to stress the system and not necessarily because they were the best algorithms possible. By way of reference, the measured speedups for **CSize-4** were 13, 11, and 8 for SOR, Matrix-Multiply, and 2D-FFT, respectively.

In general, all three programs perform best when the cluster size is four (see Figure 10). The performance of all three programs suffer when the cluster size is much larger or smaller than four. The amount of degradation reflects application-specific system service demands, due in part to the differing page-fault behavior.

At **CSize-1** we see that the SOR program degrades only slightly. The degradation is minimal because the workers have primarily independent working sets, and so there is little shared data that must be kept consistent. The increase in response time for SOR over the response time of **CSize-4** is mostly due to the increased cost of moving and destroying the worker processes on remote clusters during the initialization and termination phases.

Matrix Multiply is severely affected by too small a cluster size. This is actually a policy effect, since the default policy is to replicate read-only shared pages across clusters. Although this policy improves access times once the page is local, the cost of making 15 copies of the B matrix for **CSize-1** overshadows any gains in locality. Also, this degree of replication can cause the per-cluster working set to exceed the amount of available physical memory if the problem size is large. Using an intermediate cluster size thus helps to balance the tradeoff between locality and space utilization.

At the other end of the coupling spectrum (**CSize-16**), 2D-FFT suffers greatly because of the high number of consistency faults that occur when the workers change phases from rows to columns. Matrix Multiply is also affected by contention, because of the high number of accesses to the shared pages of the B array (see Figure 9).



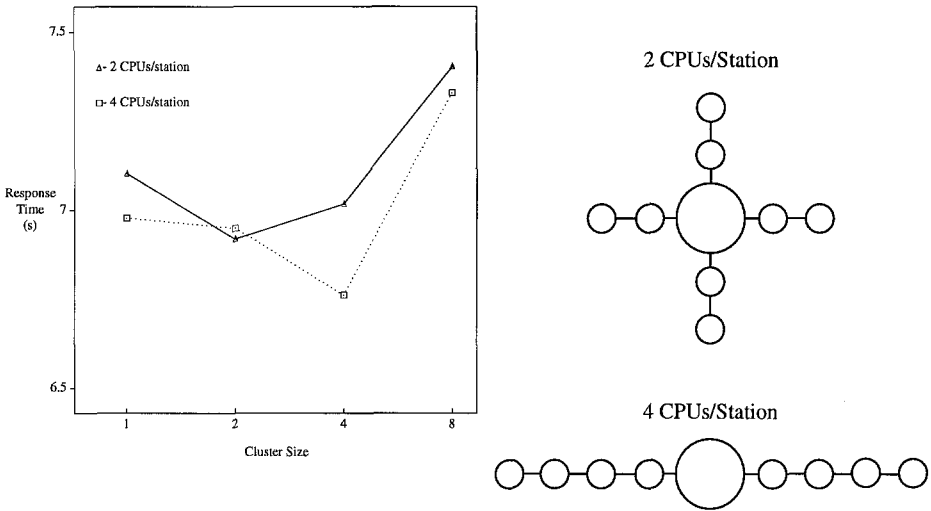


Figure 11. Matching the cluster size to the architecture for Matrix Multiply on eight processors.

Again, given the significant impact visible for a relatively small system, the benefits of clustering are likely to be more pronounced on larger systems.

**5.4. Effects of Hardware Configuration**

An interesting question is whether the cluster size should match the application communication pattern or the architectural configuration. So far, we have seen examples when both were appropriate: The synthetic stress tests of concurrent accesses to shared resources showed the cluster size should match the application; the application performance curves suggested that the cluster size should match the architecture. To explore this question further, we ran the Matrix Multiply test on an eight-processor multiprocessor configured as four hardware stations of two processors each, and configured as two hardware stations of four processors each. For both configurations the cluster size was varied from one to eight. Figure 11 shows the response times obtained for this experiment.

For both hardware configurations the figure shows that the response time is minimized when the cluster size matches the architecture. Also note that four processors per station has a generally better response time than two processors per station because there is less interstation communication. The exception occurs at a cluster size of two, since the configuration of four processors per station suffers from having two clusters sharing a common bus for independent accesses.

## 6. Lessons Learned

Our experiences with HURRICANE have demonstrated that hierarchical clustering can exploit the advantages of both small-scale tightly coupled systems and large-scale distributed systems. The tight coupling within a cluster results in performance comparable to small-scale shared-memory systems for the common case, when interactions occur primarily between objects located in the same cluster. On the other hand, because system services and data are replicated and moved to different clusters in the system, there is much locality and less sharing and hence fewer bottlenecks for independent operations. We therefore expect a system based on hierarchical clustering to scale well.

In the course of our work on HURRICANE, we have learned a number of lessons that suggest changes that could improve the performance and reduce the complexity of an operating system implemented using hierarchical clustering.

### *Locking protocols*

The locking protocols for within and across clusters were developed independently [Unrau et al. 1994], adding much complexity to the implementation. Within a cluster we enforce an ordering on the acquisition of locks to avoid deadlock. Because lock contention is often expected to be low with reasonable cluster sizes, many locks are implemented as spin locks, in which a process requesting a lock does not release the processor it is spinning on (essentially holding a lock on the processor). For intercluster operations we release all local locks, including the processor, before going remote in order to avoid deadlock.

We found that releasing locks for remote operations results in an excessively large discrepancy between the cost of local and remote operations and adds much complexity to the code. As an example of the complexity, many operations in the kernel have three versions: one that implements the local case, another for the client side of the RPC, and a third for the corresponding server side. Releasing locks for remote operations also makes it more difficult to ensure fairness, since local requests have a higher probability of acquiring a contended resource. To address the cost and fairness problems that arise from releasing locally held locks for remote operations, we modified the implementation to release locally held locks only if the locks required at the remote site are not free. This optimization resulted in even greater code complexity. We are currently investigating how to get the same benefits without unduly complicating the code.

### *Relationship between clustering and policies*

In our original implementation the cluster size determined both the structures used to manage resources and a framework for invoking policies on those resources. We have since experimented with decoupling the cluster size and the policies. For example, it is now possible to request that pages of a memory region be replicated or moved to the local processor, even if there is already a copy of the page on another processor in the

local cluster [Gamsa 1992]. For some applications we have found that this approach results in substantial performance improvements.

### *Fixed cluster sizes*

When HURRICANE is booted, a fixed-sized cluster is established for all resources. We believe that more flexibility could result in performance improvements; in particular, we believe that each class of physical (e.g., pages, processors, disks) and virtual (e.g., address spaces, regions, files) resource could be clustered independently with little increase in space requirements or complexity. As an example, the HURRICANE PPC facility manages resources more naturally on a per-processor rather than a per-cluster basis.

### *Other experiences*

In our current implementation the RPC service load is balanced by having processor  $i$  of each cluster direct its RPC requests to processor  $i$  of the target cluster. We would like to investigate more dynamic load balancing approaches; an example would be to direct an RPC request to any idle processor in the target cluster.

Some operating system services can affect all clusters on which that application is running. For example, when the state of a page used by all processes is changed, all processors must be notified. To support these operations efficiently, we have implemented an asynchronous RPC facility that allows a particular cluster to simultaneously invoke code on a number of or all other clusters. For larger-scale systems we expect to use the clustering hierarchy to implement a spanning tree broadcast that will allow not only the servicing but also the distribution of the asynchronous RPC operations to be done concurrently.

As discussed in the implementation section, local structures (representatives) are used to represent remote resources, such as processes or physical pages. In the worst case, the number of remote resources accessed by a cluster will grow proportionally to the number of clusters in the system. If we allowed the number of representatives to grow at this rate, then clearly our approach would not be scalable. However, for well-behaved applications we expect that the average number of cross-cluster interactions that each cluster will have to support will remain constant as the system size is scaled. Therefore, we believe we can obtain good performance with a fixed number of representatives per cluster, paying a somewhat higher cost to support the case when the fixed set is depleted. A proper study of the impact of this decision is still required.

## **7. Concluding Remarks**

We have introduced the concept of hierarchical clustering as a way to structure operating systems for scalability, and described how we applied this structuring technique to the HURRICANE operating system. We presented performance results that demonstrate the

characteristics and behavior of the clustered system. In particular, we showed how clustering trades off the efficiencies of tight coupling for the advantages of replication, increased locality, and reduced lock contention, but at the cost of cross-cluster overhead.

Concepts similar to hierarchical clustering have been proposed by a number of other researchers [Ahmad and Ghafoor 1991; Cheriton et al. 1991; Feitelson and Rudolph 1990; Zhou and Brecht 1991]. One of the major contributions of our work is that we have implemented a complete operating system based on this structuring methodology. To the best of our knowledge, HURRICANE is the only complete and running implementation of such a system. Our work gives a degree of confidence that systems structured in such a fashion can be reasonably implemented, a confidence that cannot be obtained from a paper design. The lessons learned from this implementation were not obvious (except perhaps in hindsight) and should influence future systems designed with similar structures.

We are now embarking on the design of a new operating system, TORNADO, for a new shared-memory multiprocessor. TORNADO will incorporate the lessons we learned from our experiences with hierarchical clustering on HURRICANE. In particular, 1) the locking protocol will be independent of the clustering, 2) each virtual and physical resource class will use its own cluster size, and 3) the cluster size and the policies will be largely decoupled.

Our experiences with scientific applications on HURRICANE/HECTOR has inspired two new research directions for TORNADO. First, we are investigating a new methodology called *physical resource-based virtual machines*, whose goal is to guarantee predictable performance, eliminating interference from other concurrently running but independent applications. This will allow compilers and applications to make optimizations that depend on predictable physical resource allocation. Second, we are investigating operating system policies and structures that will work well for applications with large data sets (e.g., data sets that do not fit in main memory).

## Acknowledgments

The research described in this paper has been partially funded by the National Sciences and Engineering Research Council of Canada and by the Information Technology Research Center of Ontario. We also wish to thank David Blythe, Yonatan Hanna, and Songnian Zhou for their significant contributions to the design and implementation of HURRICANE, and Ron White and Peter Pereira for keeping the hardware running. We also gratefully acknowledge the help of Tim Brecht, Jan Medved, Fernando Nasser, and Umakanthan Shunmuganathan.

## References

- Ahmad, I., and Ghafoor, A. 1991. Semi-distributed load balancing for massively parallel multicomputer systems. *IEEE Transactions on Software Engineering*, 17, 10 (Oct.): 987–1004.
- Anderson, T.E. 1990. The performance of spin lock alternatives for shared-memory multiprocessors. *IEEE Transactions on Parallel and Distributed Systems*, 1,1 (Jan.): 6–16.

- Balan, R., and Gollhardt, K. 1992. A scalable implementation of virtual memory HAT layer for shared memory multiprocessor machines. In *Proc., USENIX Summer '92 Conference* (San Antonio, Tex., June), pp. 107–115.
- Barach, D., Wells, R., and Uban, T. 1990. Design of parallel virtual memory management on the TC2000. Technical Report 7296, BBN Advanced Computers Inc., Cambridge, Mass.
- Barak, A., and Kornatzky, Y. 1987. Design principles of operating systems for large scale multicomputers. Technical Report RC 13220 (#59114), IBM T.J. Watson Research Center.
- BBN. 1988. *Overview of the Butterfly GP1000*. BBN Advanced Computers, Inc.
- BBN. 1989. *TC2000 Technical Product Summary*. BBN Advanced Computers, Inc.
- Bolosky, W.J., Fitzgerald, R.P., and Scott, M.L. 1989. Simple but effective techniques for NUMA memory management. In *Proc., 12th ACM Symposium on Operating System Principles*, pp. 19–31.
- Brecht, T.B. 1993. On the importance of parallel application placement in NUMA multiprocessors. In *Proc., SEDMS IV, Symposium on Experiences with Distributed and Multiprocessor Systems*, USENIX Association, pp. 1–18.
- Campbell, M., Holt, R., and Slice, J. 1991. Lock granularity tuning mechanisms in SVR4/MP. In *Proc., SEDMS II, Symposium on Experiences with Distributed and Multiprocessor Systems*, USENIX Association, pp. 221–228.
- Chaiken, D., Kubiawicz, J., and Agarwal, A. 1991. LimitLESS directories: A scalable cache coherence scheme. In *Proc., 4th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)* (Santa Clara), ACM Press, pp. 224–234.
- Chang, H.H.Y., and Rosenburg, B. 1992. Experience porting Mach to the RP3 large-scale shared-memory multiprocessor. *Future Generation Computer Systems*, 7, 2/3 (Apr.): 259–267.
- Chaves, E., Das, P.C., LeBlanc T.J., Marsh, B.D., and Scott, M.L. 1993. Kernel-kernel communication in a shared-memory multiprocessor. *Concurrency: Practice and Experience*, 5, 3, (May): pp. 171–191.
- Chen, J.B., and Bershad, B.N. 1993. The impact of operating system structure on memory system performance. In *Proc., 14th ACM Symposium on Operating Systems Principles*, pp. 120–133.
- Cheriton, D.R. 1988. The V distributed system. *Communications of the ACM*, 31,3 (Mar.): 314–333.
- Cheriton, D.R., Goosen, H., and Boyle, P. 1991. ParaDiGM: A highly scalable shared-memory multi-computer architecture. *IEEE Computer*, 24, 2 (Feb.): 33–46.
- Cox, A.L., and Fowler, R.J. 1989. The implementation of a coherent memory abstraction on a NUMA multiprocessor: Experiences with PLATINUM. In *Proc., 12th ACM Symposium on Operating System Principles*, pp. 32–44.
- Feitelson, D.G., and Rudolph, L. 1990. Distributed hierarchical control for parallel processing. *IEEE Computer*, 23, 5 (May): 65–81.
- Frank, S., Rothnie, J., and Burkhardt, H. 1993. The KSRI: Bridging the gap between shared memory and MPPs. In *IEEE Comcon 1993 Digest of Papers*, pp. 285–294.
- Gamsa, B. 1992. Region-oriented main memory management in shared-memory NUMA multiprocessors. Master's thesis, Department of Computer Science, University of Toronto, Toronto, Canada.
- Gamsa, B., Krieger, O., and Stumm, M. 1993. Optimizing IPC performance for shared-memory multiprocessors. Technical Report 294, CSRI, University of Toronto, Toronto, Canada.
- Hagersten, E., Landin, A., and Haridi, S. 1992. DDM – A cache-only memory architecture. *IEEE Computer*, 25, 9 (Sept.): 44–54.
- Krieger, O. 1994. HFS: A flexible file system for shared memory multiprocessors. Ph.D. thesis, Department of Electrical and Computer Engineering, University of Toronto, Toronto, Canada.
- Krieger, O., and Stumm, M. 1993. HFS: A flexible file system for large-scale multiprocessors. In *Proc., 1993 DAGS/PC Symposium* (Hanover, N.H., June), Dartmouth Institute for Advanced Graduate Studies, pp. 6–14.
- Krieger, O., Stumm, M., and Unrau, R. 1994. The Alloc Stream Facility: A redesign of application-level stream I/O. *IEEE Computer*, 27, 3 (Mar.): 75–83.
- Kuck, D.J., Davidson, E.S., Lawrie, D.H., and Sameh, A.H. 1986. Parallel supercomputing today, and the Cedar approach. *Science*, 231 (Feb.): 967–974.
- LaRowe Jr, R.P., Ellis, C.S., and Kaplan, L.S. 1991. Tuning NUMA memory management for applications and architectures. In *Proc., 13th ACM Symposium on Operating System Principles* (Asilomar, Pacific Grove, Calif.), Association for Computing Machinery SIGOPS, pp. 137–151.
- Lenoski, D., Laudon, J., Gharachorloo, K., Weber, W., Gupta, A., Hennessy, J., Horowitz, M., and Lam, M.S. 1992. The Stanford DASH Multiprocessor. *Computer*, 25, 3 (Mar.): 63–79.
- Mellor-Crummey, J.M., and Scott, M.L. 1991. Algorithms for scalable synchronization on shared-memory multiprocessors. *ACM Transactions on Computer Systems*, 9, 1 (Feb.): 21–65.

- Oed, W. 1993. The Cray Research massively parallel processor system CRAY T3D. Technical report, Cray Research GmbH, München, Germany.
- Peacock, J.K., Saxena, S., Thomas, T., Yang, F., and Yu, W. 1992. Experiences from multithreading system V release 4. In *Proc., SEDMS III, Symposium on Experiences with Distributed and Multiprocessor Systems*, USENIX Association, pp. 77–91.
- Pfister, G.F., Brantley, W.C., George, D.A., Harvey, S.L., Kleinfelder, W.J., McAuliffe, K.P., Melton, E.A., Norton, V.A., and Weiss, J. 1985. The IBM Research Parallel Processor Prototype. In *Proc., 1985 International Conference on Parallel Processing*, pp. 764–771.
- Scott, M.L., LeBlanc, T.J., Marsh, B.D., Becker, T.G., Dubnicki, C., Markatos, E.P., and Smithline, N.G. 1990. Implementation issues for the Psyche multiprocessor operating system. *Computing Systems*, 3, 1 (Jan.): 101–137.
- Simon, H.A. 1985. *The Sciences of the Artificial*, 2nd ed. MIT Press, Cambridge, Mass.
- Stumm, M., Unrau, R., and Krieger, O. 1992. Designing a scalable operating system for shared memory multiprocessors. In *Proc., USENIX Workshop on Microkernels and Other Kernel Architectures*, pp. 285–303.
- Unrau, R. 1993. Scalable memory management through hierarchical symmetric multiprocessing. Ph.D. thesis, Department of Electrical and Computer Engineering, University of Toronto, Toronto, Canada.
- Unrau, R., Krieger, O., Gamsa, B., and Stumm, M. 1994. Experiences with locking in a NUMA multiprocessor operating system kernel. In *Proc., USENIX OSDI Symposium* (Nov.), pp. 139–152.
- Vranesic, Z.G., Stumm, M., Lewis, D., and White, R. 1991. Hector: A hierarchically structured shared-memory multiprocessor. *IEEE Computer*, 24, 1 (Jan.): 72–80.
- Zajcew, R., Roy, P., Black, D., Peak, C., Guedes, P., Kemp, B., LoVerso, J., Leibensperger, M., Barnett, M., Rabii, F., and Netterwala, D. 1993. An OSF/1 UNIX for massively parallel multicomputers. In *Proc., USENIX Winter Conference*, USENIX Association, pp. 449–468.
- Zhou, Z., and Brecht, T. 1991. Processor pool-based scheduling for large-scale NUMA multiprocessors. In *Proc., ACM Sigmetrics Conference on Measurement and Modeling of Computer Systems* (San Diego), ACM Press, pp. 133–142.