

Nswap: A Network Swapping Module for Linux Clusters

Tia Newhall, Sean Finney, Kuzman Ganchev, and Michael Spiegel

Swarthmore College, Swarthmore, PA 10981, USA

Abstract. Cluster applications that process large amounts of data, such as parallel scientific or multimedia applications, are likely to cause swapping on individual cluster nodes. These applications will perform better on clusters with network swapping support. Network swapping allows any cluster node with over-committed memory to use idle memory of a remote node as its backing store and to “swap” its pages over the network. As the disparity between network speeds and disk speeds continues to grow, network swapping will be faster than traditional swapping to local disk. We present Nswap, a network swapping system for heterogeneous Linux clusters and networks of Linux machines. Nswap is implemented as a loadable kernel module for version 2.4 of the Linux kernel. It is a space-efficient and time-efficient implementation that transparently performs network swapping. Nswap scales to larger clusters, supports migration of remotely swapped pages, and supports dynamic growing and shrinking of Nswap cache (the amount of RAM available to store remote pages) in response to a node’s local memory needs. Results comparing Nswap running on an eight node Linux cluster with 100BaseT Ethernet interconnect and faster disk show that Nswap is comparable to swapping to local, faster disk; depending on the workload, Nswap’s performance is up to 1.7 times faster than disk to between 1.3 and 4.6 times slower than disk for most workloads. We show that with faster networking technology, Nswap will outperform swapping to disk.

1 Introduction

Using remote idle memory as backing store for networked and cluster systems is motivated by the observation that network speeds are getting faster more quickly than are disk speeds [9]. In addition, because disk speeds are limited by mechanical disk arm movements and rotational latencies, this disparity will likely grow. As a result, swapping to local disk will be slower than using remote idle memory as a “swap device” and transferring pages over the faster network. Further motivation for network swapping is supported by several studies [2,4,10] showing that large amounts of idle cluster memory are almost always available for remote swapping.

We present Nswap, a network swapping system for heterogeneous Linux clusters and networks of Linux machines. Nswap transparently provides network swapping to cluster applications. Nswap is implemented as a loadable kernel

module that is easily added as a swap device to cluster nodes and runs entirely in kernel space on an unmodified ¹ Linux kernel; applications can take advantage of network swapping without having to re-compile or link with special libraries. Nswap is designed to scale to large clusters using an approach similar to Mosix's design for scalability [1]. In Nswap there is no centralized server that chooses a remote node to which to swap. In addition, Nswap does not rely on global state, nor does it rely on complete or completely accurate information about the state of all cluster nodes to make swapping decisions. Thus, Nswap will scale to larger clusters because each node independently can make swapping decisions based on partial and not necessarily accurate information about the state of the cluster.

Nswap supports dynamic growing and shrinking of each node's Nswap cache in response to the node's local memory use, allowing a node to reclaim some of its Nswap cache space for local paging or file I/O when it needs it, and allowing Nswap to reclaim some local memory for Nswap cache when the memory is no longer needed for local processing. Growing and shrinking of the Nswap cache is done in large chunks of memory pages to better match the bursty behavior of memory usage on a node [2]. Another feature of Nswap is that it avoids issuing writes to disk when swapping; it does not write-through to disk on swap-outs, and it supports migration of remotely swapped pages between the cluster nodes that cache them in response to changes in a remote node's local memory use. Only when there is no available idle memory in the cluster for remote swapping does Nswap revert to swapping to disk. Because Nswap does not do write-through to disk, swapping activity on a node does not interfere with simultaneous file system disk I/O on the node. In addition, as long as there is idle remote memory in the cluster, Nswap can be used to provide disk-less cluster nodes the ability to swap.

In section 2 we discuss related work in network swapping. In section 3 we present the details of Nswap's implementation, including a presentation of the swapping protocols. In section 4 we present the results of measurements of Nswap running on an eight node cluster. Our results show that with 100 BaseT technology we are comparable to high-end workstation disks, but that with faster interconnect, Nswap will outperform swapping to disk. In section 5 we discuss future directions for our work.

2 Related Work

There have been several previous projects that examine using remote idle memory as backing store for nodes in networks of workstations [3,6,10,8,5,11,7]. For example, Feeley et. al.[6] implement a network cache for swapped pages by modifying the memory management system of the DEC OSF/1 kernel. Their system views remote memory as a cache of network swapped pages that also are written through to disk; remote servers only cache clean pages and can arbitrarily drop

¹ Currently, we require a re-compile of the 2.4.18 kernel to export two kernel symbols to our module, but we have not modified Linux kernel code in any way.

a page when their memory resources become scarce. Each node's memory is partitioned into local space and network cache space. When a page fault occurs, a node's local space may grow by one page at the expense of a page of its network cache space. Our growing and shrinking policies do something similar, but we grow and shrink in larger units to better match bursty memory use patterns. In addition, our system does not do a write-through to disk on every swap-out. Thus, our system will not be slowed down by disk writes and will not interfere with a node's local file I/O.

Markatos and Dramitinos [10] describe reliability schemes for a remote memory pager for the DEC OSF/1 operating system. Their system is implemented as a client block device driver, and a user-level server for storing pages from remote nodes. Both dirty and clean pages can be stored at remote servers. When a server is full, or when it needs more memory for local processes, remote pages are written to the server's disk, resulting in a significant slow down to a subsequent swap-in of the page. Our system avoids writing pages to disk in similar circumstances by migrating the page to another server. In addition, our system allows a server's cache size to change, and it allows nodes to dynamically change roles between clients and servers based on their local memory use.

Bernard and Hamma's work [5] focuses on policies for balancing a node's resource usage between remote paging servers and local processes in a network of workstations. Their policy uses local and remote paging activity and local memory use to determine when to enable or disable remote paging servers on a node. Our growing and shrinking policies for Nswap cache sizes have to solve a similar problem. However, on a cluster system, where there is no notion of a user "owning" an individual cluster node, so we expect that we can keep Nswap cache space on nodes with under-committed memory but higher cpu loads.

3 Nswap Implementation

Nswap consists of two components running entirely in kernel space (shown in Figure 1). The first is a multi-threaded client that receives swap-in and swap-out requests from the kernel. The second is a multi-threaded server that manages a node's Nswap cache and handles swap-in and swap-out requests from remote clients. Each node in the Nswap cluster runs an Nswap client and an Nswap server. At any given time, a cluster node is acting either as a client or a server, but typically not as both simultaneously; a node adjusts its role based on its local memory use. The goals of Nswap's design are to provide efficient, transparent network swapping to cluster applications, to dynamically adjust to individual node's memory needs, and to scale to large, heterogeneous clusters.

3.1 Nswap Client and Server

The Nswap client is implemented as a block pseudo-device. This abstraction allows an easy and portable interface to the Linux kernel's swapping mechanism, which communicates with Nswap exactly as it would any physical block device. Idle memory across the cluster is accessed through read and write requests from

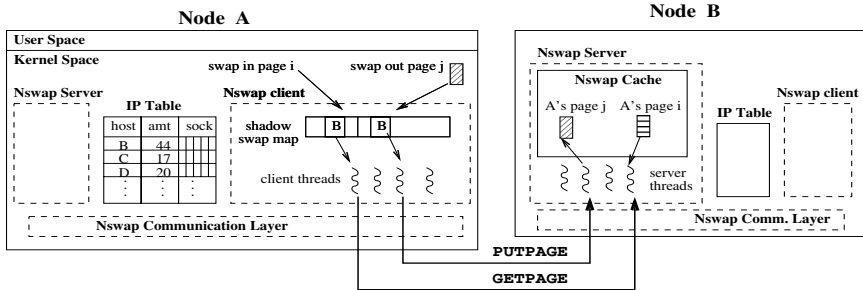


Fig. 1. Nswap System Architecture. Node A shows the details of the client including the shadow slot map used to store information about which remote servers store A's pages. Node B shows the details of the server including the Nswap cache of remotely swapped pages. In response to the kernel swapping in (out) a page to our Nswap device, a client thread issues a PUTPAGE (GETPAGE) to write (read) the page from a remote server.

the kernel to the Nswap client device. The client uses multiple kernel threads to simultaneously service several swap-in or swap-out requests.

The Nswap client needs to add additional state to the kernel's swap map so it can keep track of the location of its remotely swapped pages. We add a shadow swap map (Figure 1) that stores the following information (in 16 bits) for each slot in the kernel's swap map of our device: the **serverID**, designating which remote server stores the page; the **hop_count**, counting the number of instances a page has been migrated (used to avoid race-conditions that can occur during page migration and to limit the number of times a page can be migrated); the **time_stamp**, identifying an instance of the slot being used by the kernel (used to identify "dead" pages in the system that can be dropped by the server caching them); and the **in_use** bit to control conflicting simultaneous operations to the same swap slot. Additionally, these fields allow for a communication protocol that does not require the client's state to be synchronously updated during page migrations.

The Nswap server is implemented as a kernel-level daemon that makes local memory available for Nswap clients. When the Nswap module is loaded, the server starts three kernel threads: a listener thread, a memory thread, and a status thread. The listener thread listens for connections from Nswap clients, and starts new threads to handle the communication. The memory thread monitors the local load, communicating its load to other hosts, and if necessary, triggering growing and shrinking its Nswap cache. The status thread accepts UDP broadcast messages from memory threads on other servers. These messages contain changes in a server's available Nswap cache size. The status thread updates its IPTable with this information.

The IPTable on each node contains a potentially incomplete list of other server's state. Each entry contains an Nswap server's IP, the amount of Nswap cache it has available, and a cache of open sockets to the server so that new

connections do not have to be created on every page transfer. The Nswap client uses information in the IPTable to select a remote server to send its pages, and to obtain connections to remote servers. The information in the IPTable about each server's Nswap cache availability does not have to be accurate, nor does there need to be an IPTable entry for every node in the cluster, for a client to choose a remote server on a swap-out. This design will help Nswap scale better to larger clusters, but at the expense of clients perhaps not making the best possible server choice. Typically, the information is accurate enough to make a good server choice. Nswap recovers from a bad server choice by migrating pages to a better server.

When the client makes a request for a page, it sends page meta-data (`clientID`, `slot_number`, `hop_count`, `timestamp`) that the server uses to locate the page in its Nswap cache. Page meta-data is also used to determine when cached pages are no longer needed. Since the Linux kernel does not inform a swap device when it no longer needs a page, "dead" pages can accumulate. A garbage collector thread in the Nswap client periodically runs when the client has not swapped recently. It finds and cleans dead slots in the swap slot map, sending the server a message indicating that it can drop the "dead" page from its cache. In addition, "dead" pages are identified and dropped during page migration, and when a client re-uses an old slot during a swap-out operation.

A predictive monitoring policy is used to dynamically grow or shrink the size of Nswap cache. A status thread on each node periodically polls the behavior of the Linux memory management subsystem. If the machine is showing symptoms of high memory usage, Nswap cache will shrink, possibly triggering migration of some of the remote pages it caches to other Nswap servers. When local memory is underutilized, Nswap cache size is increased.

3.2 Nswap Communication Protocol

The communication protocol between different nodes in an Nswap cluster is defined by five types of requests: **PUTPAGE**, **GETPAGE**, **PUNTPAGE**, **UPDATE** and **INVALIDATE**. When Nswap is used as a swap device, the kernel writes pages to it. The client receives a write request and initiates a **PUTPAGE** to send the page to a remote server. At some later time the kernel may require the data on the page and issue a read request to the Nswap device. The client uses the **GETPAGE** request to retrieve the page (see Figure 1). If the workload distribution changes, it may become necessary for a server to reduce its Nswap cache size, using the **PUNTPAGE** request to offload the page to another Nswap server. Moving a page from one server to another involves an **UPDATE** request to alert the client to the new location of the page and an **INVALIDATE** request from the client to the old server to inform the old server that it can drop its copy of the page (see Figure 2).

GETPAGE and **PUTPAGE** are designed to be as fast as possible for the client who is currently swapping. If a client makes a bad server choice for a **PUTPAGE**, the Nswap servers handle it through page migration rather than forcing the client to make a better server choice, which would slow down the client's swap-outs. In addition, the protocol is designed to limit synchronous activities. As a result,

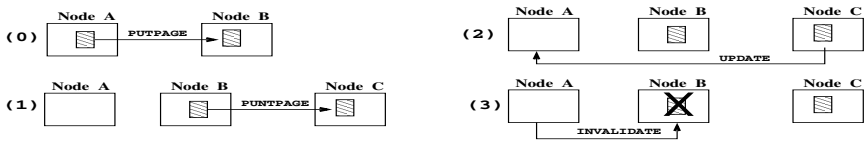


Fig. 2. Page Migration in Nswap. Node A acts as an Nswap client, and Nodes B and C act as servers. A **PUTPAGE** from server B to server C (1), triggers an **UPDATE** from server C to client A (2), which in turn triggers an **INVALIDATE** from client A to server C (3). At this point, B can drop its copy of A's page.

extra state (page meta-data) is passed with requests so that a receiver can detect out-of-order and old requests and handle them appropriately. The following is an overview of the communication protocol:

PUTPAGE is used to ask a remote host to store a local page. The client picks a remote server using its IPTable information, and sends the server a **PUTPAGE** command and the page's meta-data. The server almost always responds with an **OK_PUTPAGE**, and then the client sends the page data. Even a server with a full Nswap cache typically will accept the page from the client. In these cases, the server starts to remove some of its least recently cached pages by issuing **PUNTPAGES** after completing the **PUTPAGE** transaction.

GETPAGE is a request to retrieve a page that is stored remotely. The client sends the **GETPAGE** command and the page meta-data. The server uses the meta-data to find the page in its Nswap cache, and sends it to the client.

INVALIDATE is used by a client to inform a server that it may drop a page from its Nswap cache. An **INVALIDATE** can result from page migration (see **PUNTPAGE** below) or from garbage collection on the client. The client sends an **INVALIDATE** command and the page meta-data. The server compares all the meta-data of the page and frees the page if these match. Otherwise, the **INVALIDATE** request is for an old page and is ignored.

PUNTPAGE is used to implement page migration. When a server becomes over-committed, it attempts to get rid of the least recently cached foreign pages by punting them to other servers. The over-committed server sends a **PUNTPAGE** command and the page data to another remote server. The original server cannot drop the page until it receives an **INVALIDATE** from the client as described above. Once the page has been transferred to the new server, the new server initiates an **UPDATE** to the owner of the page. If there is no available Nswap cache space in the cluster, the server punts the page back to its owner who writes it to its local disk.

UPDATE is used to inform a client that one of its pages has moved. The new server sends an **UPDATE** command to the client with the page meta-data. The page meta-data is used to detect whether the **UPDATE** is for an old page, in which case the client sends an **INVALIDATE** to the sender of the **UPDATE**; otherwise, the client sends an **INVALIDATE** to the previous server caching the page so that it can drop its copy of the page.

4 Results

We present results comparing swapping to disk and Nswap for several workloads. Our experiments were run on a cluster of 8 nodes² running version 2.4.18 of the Linux kernel connected by 100BaseT Ethernet. Four of the nodes have Intel Pentium II processors, 128 MB of RAM, and Maxtor DiamondMax 4320 disks with a data rate of up to 176 Mb/sec. The other machines have Pentium III processors, 512 MB of RAM, and IBM Deskstar disk with a sustained data rate of 167-326 Mb/sec and a max rate of 494 Mb/sec. In both machine types, disk transfer rates are faster than network transfer rates, so we expect to be slower than swapping to disk. However, our results show that for several workloads swapping across a slower network is faster than swapping to faster local disk. In addition, we calculate that on 1 and 10 Gbit Ethernet, Nswap will outperform swapping to disk for almost all workloads.

4.1 Workload Results

Table 1 shows run times of several workloads comparing Nswap with swapping to disk. On the PIII's we ran two versions of Nswap, one using TCP/IP to transfer pages (column 5), the other using UDP/IP (column 6). We ran one and four process versions of each workload with one node acting as the client and three nodes acting as servers. We also disabled Nswap cache growing and shrinking to reduce the amount of variation between timed runs.

Workload 1 consists of a process that performs a large sequential write to memory followed by a large sequential read. It is designed to be the best case for swapping to disk because there will be a minimal amount of disk head movement when swapping due to the way in which Linux allocates space on the swap partition. Workload 2 consists of a process that performs random writes followed by random reads to a large chunk of memory. This workload stresses disk head movement within the swap partition only. Workload 3 consists of two processes. The first runs a Workload 1 application and the second is a process that performs a large sequential write to a file. Workload 3 further stresses disk head movement when swapping to disk because file I/O and swap I/O are concurrently taking place in different disk partitions. Workload 4 consists of two processes; the first runs the Workload 2 application and the second runs the large sequential file write application. The four process versions of each workload are designed to represent a more realistic cluster workload. The size of each Workload varied from platform to platform based on the size of RAM on each machine, and they varied between the Random and Sequential tests. As a result, for each Workload (row), only the values in columns 2 and 3 can be compared, and the values in columns 4, 5, and 6 can be compared.

We expect that disk will perform much better than Nswap on the single process versions of Workload 1 and Workload 2 because the disk arm only moves

² Although our cluster currently consists of only Pentium architectures, Nswap will run on any architecture supported by Debian Linux.

Table 1. Swapping to fast disk vs. TCP Nswap and UDP Nswap on 100BaseT for PII and PIII nodes. The rows are 1 and 4 process runs of each of the four Workloads. Time is measured in seconds and the values are the mean of 5 runs of each benchmark. Bold values indicate runs for which Nswap is faster than disk.

(Workload) # procs	DISK (PII)	NSWAP (PII, TCP)	DISK (PIII)	NSWAP (PIII, TCP)	NSWAP (PIII, UDP)
(1), 1	98.0	450.9 (4.6x slower)	13.1	154.3 (11.7x slow)	61.3 (4.6x slow)
(1), 4	652.9	630.6 (1.04x fast)	551.4	1429.7 (2.6x slow)	614.4 (1.1x slow)
(2), 1	874.6	1937.0 (2.2x slow)	266.8	1071.8 (4.0x slow)	153.5 (1.7x fast)
(2), 4	996.7	617.0 (1.6x fast)	68.6	189.3 (2.8x slow)	50.3 (1.4x fast)
(3), 1	632.9	737.7 (1.7x slow)	770.2	1111.1 (1.4x slow)	811.0 (1.1x slow)
(3), 4	1312.1	1127.2 (1.2x fast)	727.1	1430.5 (1.9x slow)	619.5 (1.2x fast)
(4), 1	1971.2	2111.0 (1.07x slow)	923.9	1529.3 (1.7x slow)	821.7 (1.1x fast)
(4), 4	1453.0	1094.6 (1.3x fast)	502.5	498.7 (1.01x fast)	429.2 (1.2x fast)

within the swap partition when these workloads run, and our results confirm this; the worst slow downs (4.6 and 11.7) are for the single process runs of Workload 1. The differences between the performance of Workloads 1 and 2 show how disk head movement within the same disk partition affects the performance of swapping to disk. For the four process runs of Workloads 1 and 2, Nswap is much closer in performance to disk, and is faster in some cases; the TCP version on the PIIIs is 2.6 and 2.8 times slower than disk, the UDP version on the PIIIs is 1.4 times faster for Workload 2, and on the PIIIs both workloads are faster than disk (1.04 and 1.6 times faster). This is due to a slow down in disk swapping caused by an increase in disk arm movement because multiple processes are simultaneously swapping to disk. The results from Workloads 3 and 4, where more than one process is running and where there is concurrent file I/O with swapping, show further potential advantages of Nswap; Nswap is faster than disk for the four process version of all Workloads on the PIIIs, for three of the four Workloads under UDP Nswap, and for Workload 4 under TCP Nswap on the PIIIs. The differences in TCP and UDP Nswap results indicate that TCP latency is preventing us from getting better network swapping performance. Nswap performs better for Workloads 3 and 4 because network swapping doesn't interfere with file system I/O, and because there is no increase in per-swap overhead when multiple process are simultaneously swapping.

4.2 Results for Nswap on Faster Networks

Currently we do not have access to faster network technology than 100BaseT Ethernet, so we are unable to run experiments of Nswap on a cluster with a network that is actually faster than our cluster's disks. However, based on measurements of our workloads running Nswap with 10BaseT and Nswap with 100BaseT Ethernet, we estimate run times on faster networks.

The rows of Table 2 show execution times of workloads run on 100BaseT and 10BaseT, and our estimates of their run times on 1 and 10 Gigabit Ethernet.

Table 2. Time (and Speedups) for TCP & UDP Nswap on faster networks. All Workloads were run on the PIIIs. Bold values indicate cases when Nswap is faster than disk. The rows show for each workload calculated 1 and 10 Gbit run times (col. 5 & 6) based on measured speedup values between runs on 10 and 100 BaseT (col. 3 & 4).

(Workload)	Disk	10BaseT	100BaseT	1 Gbit	10 Gbit
(1) TCP	580.10	5719.00	1518.34 (speedup 3.8)	1075.00 (5.3)	1034.17 (5.5)
(1) UDP	12.27	306.69	56.80 (speedup 5.4)	28.90 (10.6)	26.30 (11.6)
(2) UDP	226.79	847.74	153.54 (speedup 5.5)	77.30 (10.9)	70.30 (12.1)
(4) UDP	6265.39	9605.91	1733.93 (speedup 5.54)	866.18 (11.1)	786.72 (12.2)

We use the speedup results from our 10 and 100 BaseT measurements of each benchmark and apply Amdahl’s Law to get estimates of speedups for 1 and 10 Gigabit Ethernet using the following equation:

$$TotalSpeedup = \frac{1}{1 - FractionBandwith + FractionBandwith / SpeedupBandwidth}.$$

For each Workload, we compute *FractionBandwith* based on *TotalSpeedup* values from measured 10 and 100 BaseT runs. Using this value, we compute *TotalSpeedup* values for 1 and 10 Gbit Ethernet (columns 5 and 6 in Table 2).

The measured speedups between 10 and 100 BaseT of TCP version of Nswap are lower than we expected (one example is shown in row 1 of Table 2); timed runs of these workloads when they completely fit into memory compared to when they don’t, show that over 99% of their total execution time is is due to swapping overhead. We also measured our `GETPAGE` and `PUTPAGE` implementations and found that over 99% of their time is due to transferring page data and metadata. However, the UDP version of Nswap results in speedup values closer to what we expected. In fact, UDP Nswap on Gigabit Ethernet outperforms swapping to disk for all Workloads except the single process version of Workload 1 on the PIIIs. Our speedup results further indicate that TCP latency is preventing Nswap from taking full advantage of improvements in network bandwidth.

5 Conclusions and Future Work

Nswap is a network swapping system for heterogeneous Linux clusters. Because Nswap is implemented as a loadable kernel module that runs entirely in kernel space, it efficiently and transparently provides network swapping to cluster applications. Results from experiments on our initial implementation of Nswap on an eight node cluster with 100BaseT Ethernet show that Nswap is comparable to swapping to faster disk. Furthermore, since it is likely that network technology will continue to get faster more quickly than disk technology, Nswap will be an even better alternative to disk swapping in the future.

One area of future work involves examining reliability schemes for Nswap. We are also investigating reliable UDP implementations for the Nswap communication layer so that we can avoid the latency of TCP to better take advantage

of faster networks, but to still get reliable page transfers over the network. Other areas for future work involve testing Nswap on faster networks and larger clusters, and further investigating predictive schemes for determining when to grow or shrink Nswap caches sizes. In addition, we may want to examine implementing an adaptive scheme into Nswap. Based on the results from the single process version of Workload 1 in Table 1, there may be some workloads for which swapping to disk will be faster. Nswap could be designed to identify these cases, and switch to disk swapping while the workload favors it.

Acknowledgments. We thank Matti Klock, Gabriel Rosenkoetter, and Rafael Hinojosa for their participation in Nswap's early development.

References

1. Barak A., La'adan O., and Shiloh A. Scalable cluster computing with MOSIX for Linux. In *Proceedings of Linux Expo '99*, pages 95–100, Raleigh, N.C., May 1999.
2. Anurag Acharya and Sanjeev Setia. Availability and Utility of Idle Memory on Workstation Clusters. In *ACM SIGMETRICS Conference on Measuring and Modeling of Computer Systems*, pages 35–46, May 1999.
3. T. Anderson, D. E. Culler, D. A. Patterson, and the NOW Team. A case for NOW (Networks of Workstations). *IEEE Micro*, February 1999.
4. Remzi H. Arpaci, Andrea C. Dusseau, Amin M. Vahdat, Lok T. Liu, Thomas E. Anderson, and David A. Patterson. The Interaction of Parallel and Sequential Workloads on a Network of Workstations. In *ACM SIGMETRICS Conference on Measurement and Modeling of Computer Systems*, pages 267–278, 1995.
5. G. Bernard and S. Hama. Remote Memory Paging in Networks of Workstations. In *SUUG'94 Conference*, April 1994.
6. Michael J. Feeley, William E. Morgan, Frederic H. Pighin, Anna R. Karlin, Henry M. Levy, and Chandramohan A. Thekkath. Implementing Global Memory Management in a Workstation Cluster. In *15th ACM Symposium on Operating Systems Principles*, December 1995.
7. Michail D. Flouris and Evangelos P. Markatos. *Network RAM, in High Performance Cluster Computing: Architectures and Systems, Chapt. 16*. Prentice Hall, 1999.
8. Liviu Iftode, Karin Petersen, and Kai Li. Memory Servers for Multicomputers. In *IEEE COMPCON'93 Conference*, February 1993.
9. John L. Hennessy and David A. Patterson. *Computer Architectures A Quantitative Approach, 3rd Edition*. Morgan Kaufman, 2002.
10. Evangelos P. Markatos and George Dramitinos. Implementation of a Reliable Remote Memory Pager. In *USENIX 1996 Annual Technical Conference*, 1996.
11. Li Xiao, Xiaodong Zhang, and Stefan A. Kubricht. Incorporating Job Migration and Network RAM to Share Cluster Memory Resources. In *Ninth IEEE International Symposium on High Performance Distributed Computing (HPDC'00)*, 2000.