

## *MPVM: A Migration Transparent Version of PVM*

Jeremy Casas, Dan L. Clark, Ravi Konuru,  
Steve W. Otto, Robert M. Prouty, and  
Jonathan Walpole  
Oregon Graduate Institute of Science & Technology

---

**ABSTRACT:** Parallel Virtual Machine (PVM) is a widely-used software system that allows a heterogeneous set of parallel and serial UNIX-based computers to be programmed as a single message-passing parallel machine. In this paper, an extension to PVM to support dynamic process migration is presented. Support for migration is important in general-purpose workstation environments since it allows parallel computations to co-exist with other applications, using idle-cycles as they become available and off-loading from workstations when they are no longer free. A description and evaluation of the design and implementation of the prototype Migratable PVM system is presented together with some performance results.

---

## *1. Introduction*

PVM [Beguelin et al. 1993.7; Dongarra et al. 1993; Beguelin et al. 1993.6] is a software system that allows a heterogeneous network of parallel and serial computers to be programmed as a single computational resource. This resource appears to the application programmer as a potentially large message-passing virtual computer. Such a system allows the computing power of widely available, general-purpose computer networks to be harnessed for parallel processing. With the rapid advances in workstation performance, such networks already provide a viable and affordable alternative to expensive special-purpose supercomputers.

General-purpose workstation networks have certain key characteristics that must be considered when they are to be used for parallel processing. First, the collective resources of the network are often shared by a potentially large number of users running a wide range of applications. Second, despite the high level of sharing, the concept of ownership is frequently present. In particular, individual workstations, while available across the network, are likely to be owned by some specific user or have priority users. Workstation owners are often willing to allow others to access their workstation when it is idle, but expect dedicated access the rest of the time. Since most workstations are idle most of the time [Litzkow et al. 1988], the key to harnessing the full power of such systems lies in gaining access to these idle cycles.

For PVM to gain unobtrusive access to idle cycles, it must be able to (a) recognize when a workstation becomes available for it to use, (b) recognize when a workstation ceases to be available to it, and (c) migrate processes between workstations so that work can be assigned to newly available workstations and off-loaded from workstations that are being reclaimed by their owners. Automatic and timely off-loading of processes requires PVM to be extended to support dynamic process migration.

In this paper, Migratable PVM (MPVM), an extension of PVM that allows parts of the parallel computation to be suspended and subsequently resumed on other workstations, is presented. There were three key goals under consideration in the design of MPVM. First, migration had to be transparent to both application programmer and user. Neither the programmer nor the user needs to know that

portions of the application are migrating. Second, source-code compatibility with PVM had to be maintained. Source-code compatibility would allow existing PVM applications to run under MPVM without, or at least with minimal, modification. Lastly, MPVM had to be as portable as possible.

The remainder of the paper is organized as follows. Section 2 gives an overview of PVM and the problem addressed by this work. Section 3 outlines the design and implementation of MPVM, and is followed by performance results in Section 4. Related work is presented in Section 5, a qualitative discussion of the design and implementation in Section 6, and conclusions and future work in Section 7.

## *2. Background*

MPVM is based on PVM 3.3.4 as released from Oak Ridge National Laboratory and is part of the larger Concurrent Processing Environment, an ongoing research effort [Beguelin et al. 1993.6]. This section presents an overview of the PVM system and the problems that had to be addressed to support task migration.

### *2.1. PVM Overview*

The PVM system consists of a daemon process called the `pvmd` running on each host on a network of workstations and a run-time library called the `pvmlib` linked into each application process (Figure 1). Each `pvmd` is assigned a unique host ID or `hid`. The `pvmlib` defines a suite of PVM primitives that presents a “message-passing parallel machine” user-interface to the application.

A PVM application is composed of Unix processes linked with the `pvmlib`. These processes, called tasks in PVM, communicate with each other via message-passing primitives found in the `pvmlib`. Just like the `pvmds`, each task is assigned a task ID or `tid` which uniquely identifies each task in the virtual machine. These `tids` are used to designate the source and destination tasks for messages (i.e., messages are addressed to tasks, not to ports or mailboxes).

Messages passed within the PVM system can be categorized into system messages and application messages. System messages are used exclusively by PVM to manage the virtual machine and perform application code requests (e.g., spawn a new task, get information about the virtual machine, etc.). The application code is not aware of these messages. Application messages on the other hand are used exclusively by the PVM application.

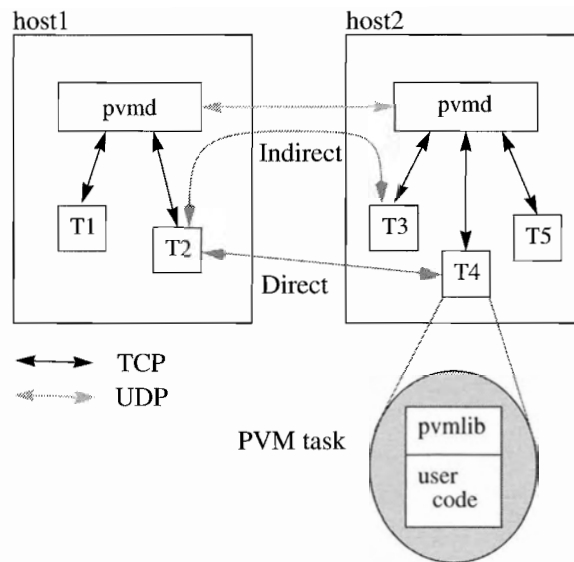


Figure 1. The PVM system. PVM is composed of daemons (*pvmds*) running on each host of the virtual machine and a run-time library (*pvmlib*) linked into every task.

PVM provides two routing mechanisms for application messages: indirect and direct routing. The choice of routing mechanism to use is controlled by the application code. By default, messages are routed indirectly. Using indirect routing, as illustrated in Figure 1, a message from task T2 to T3 passes through T2's local *pvmd* (*pvmd* on host1), through T3's local *pvmd* (*pvmd* on host2), and finally to T3. *Pvmd*-to-*pvmd* communication uses UDP socket connections while task-to-*pvmd* communications use a TCP socket connection which is established during task start-up. In direct routing, a message from task T2 to T4, also illustrated in Figure 1, uses a TCP socket connection between T2 and T4, by-passing the *pvmds* altogether. TCP connections between tasks are created "on-demand". Only when tasks that have set their routing option to use direct routing start communicating with each other are TCP connections established.

An important aspect to remember when using PVM is the message-ordering semantics it provides. PVM guarantees that messages sent from one task to another are received in the same order they were sent. The importance of recognizing this "guarantee" is that there are PVM applications that take advantage of this message-ordering semantics. Hence, new versions of PVM such as MPVM should maintain the same semantics.

Lastly, in PVM 3.3.4, it is possible to designate a special task as the resource manager. The resource manager, also called the global scheduler (GS) in this

paper, embodies decision making policies [Al-Saqabi et al. 1994] such as task-to-processor allocation for sensibly scheduling multiple parallel applications. Using a GS makes it convenient to experiment with different scheduling policies. In MPVM, the interface between the pvmds and the GS has been extended to accommodate task migration, allowing the GS to use dynamic scheduling policies.

## 2.2. PVM Task Migration: the Problem

Task migration is the ability to suspend the execution of a task on one machine and subsequently resume its execution on another. A major requirement for task migration is that the migration should not affect the correctness of the task. Execution of the task should proceed as if the migration never took place. To ensure the “transparency” of the migration, it is necessary to capture the state of the task on the source machine and reconstruct it on the target machine.

The state of a task can be viewed in two ways: its state as a Unix process and its state as a task of a PVM application. From the point of view of the operating system (OS), a task is just a single process. As such, its state includes (1) the processor state, (2) the state held by the process, (3) the state held by the OS *for* the process, and (4) the state held *by* the process about the local OS.

The processor state includes the contents of the machine registers, program counter, program status word, etc. This information defines exactly where the task was executing prior to migration, and consequently, where execution should resume upon restart on the target machine. The state held by the process itself includes the contents of its text, data (static and dynamic), and stack. The state held by the OS for the process includes signal information (e.g., blocked signals, pending signals), open files, and socket connections to name a few. Other less obvious state information held by the OS includes page table entries, controlling terminals, and process relationship information (e.g., parent/child process relationship and process groups). OS state held by the process includes file descriptors, process IDs, host name, and time. These are state information, known to the process, that are only valid in the context of the local execution environment (local OS and host).

From the point-of-view of PVM, a task is one of a set of tasks that makes up an application. In this context, a task’s state includes its tid and the messages sent to/from that task. Regardless of migration, each task should be referred to using the same tid, no message should be lost, and all messages should be received in the correct order (as defined by PVM).

Thus, the problem addressed by MPVM is how to capture and reconstruct the state information so that tasks can be migrated from one machine to another without affecting the correctness of the entire application.

### 3. Design and Implementation

In this section, the design and implementation of MPVM is described. In order to support task migration, both the pvmd and pvmlib had to be modified. The modifications made were also driven by the goals of source code compatibility, portability, and migration transparency. To ensure source code compatibility, the modifications had to maintain the same function calls, parameters and semantics, as provided by PVM. To maximize portability, the migration mechanism had to be implemented at user-level, using facilities available through standard UNIX library routines and system calls. Migration transparency is addressed by modifying the pvmd and pvmlib such that the migration could occur without notifying the application code and by providing “wrapper” functions to certain system calls. A more complete evaluation of these goals are presented in Section 6.

#### 3.1. Application Start-up

The primary interface to the migration mechanism in MPVM is through the signal mechanism provided by UNIX. That is, task migration is initiated using a migration signal sent from the pvmd to the migrating task. The migrating task should have a migration signal handler installed to catch the migration signal. At this point, it is only important to know that a signal handler has to be installed for migration to work. The function of the migration signal handler will be discussed in Section 3.2.

To avoid explicitly modifying the source code of the PVM application to install the signal handler, the pvmlib defines its own `main()` function which executes the necessary initialization and then calls a function called `Main()`. When the application program is compiled, its `main()` function is “renamed” to `Main()` using “C” macro substitution facilities available through the compiler (e.g., `-Dmain=Main` flag). Thus, when the application code is linked with the pvmlib, the resulting executable will have the pvmlib’s `main()` as the entry point, allowing execution of the migration initialization code prior to the execution of the application’s code.

While this solution is simple, not to mention inherently portable, it will fail when combined with other systems that use the same “trick” for executing code prior to the application’s `main()`. An alternative solution is to define a customized version of the start-up code, usually `crt0.o` (“C” Run-Time object module).

### 3.2. Migration Protocol

Once the application is up and running, it executes just as a traditional PVM application would, until a task has to migrate. There are a number of reasons for a task to migrate: excessively high machine load, machine reclamation by its owner, a more suitable machine becomes available, etc. Regardless of the rationale for migration, the same migration mechanism can be used.

A migration protocol is used to facilitate the migration. The migration protocol is divided into four stages as shown in Figure 2. While the first stage addresses “when” migration occurs, the last three stages correspond exactly to the main components of migration: state capture, transfer, and reconstruction.

An important component of the migration protocol is what is collectively called Control Messages. These control messages or CMs are special system messages added to the pvmds and the pvmlib for the primary purpose of managing task migration. Just like other system messages, these control messages are

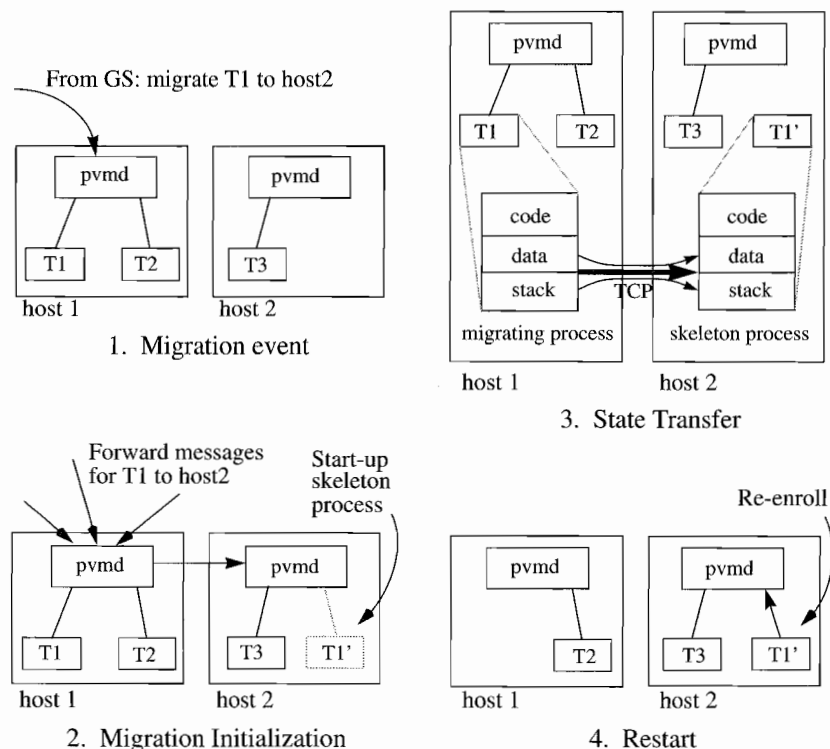


Figure 2. Migration protocol. This figure illustrates the stages involved in migrating task T1 from host1 to host2.

invisible to the application code. There are different kinds of CMs, each of which will be discussed in the following sections.

### 3.2.1. Migration Event

The migration of a task is triggered by a migration event. This event triggers the GS which determines whether or not tasks have to be migrated. If so, it also decides which tasks to migrate and to where.

If the GS decides to migrate a task, an *SM\_MIG* CM (SM stands for Scheduler Message) is sent by the GS to the pvmd on the host where the task to be migrated is currently executing. This *SM\_MIG* CM contains a tid and an hid, indicating the task to be migrated and the destination host respectively. For brevity, the task to be migrated shall be referred to as Mtask, the pvmd on the host where Mtask will be migrating from as Spvmd, and the pvmd on the destination host as Dpvmd.

### 3.2.2. Migration Initialization

Upon receipt of an *SM\_MIG* CM, the Spvmd verifies that the tid belongs to a locally executing task and that the hid refers to a valid host (not itself). If either of the tid/hid is invalid, a PvmNoTask/PvmNoHost error code is sent back to the GS via an *SM\_MIGACK* CM.

Migration initialization is divided into two components which occur in parallel. The first component, local initialization, involves “priming-up” Mtask for the state transfer. The second component, remote initialization, involves the creation of a “skeleton process” that will be the recipient of the state information to be transferred (Figure 3).

Local initialization begins when a SIGMIGRATE signal is sent from the Spvmd to Mtask (step 1' in Figure 3). This signal is caught by the migration signal handler installed by Mtask during its start-up (recall the application start-up discussion). The advantage of using a migration signal handler is two-fold: first, it allows for asynchronous task migration, and second, it is the main component used in capturing the processor state. When the signal handler is invoked, the OS automatically stores a copy of the processor state at the time when the process was interrupted in the user stack. This saved processor state is used to restore the state of execution of the process/task upon return from the signal handler. Currently, SIGTSTP is used as the SIGMIGRATE signal.

To prepare for the process state transfer in the next stage, the migration signal handler in Mtask flushes all messages in the TCP socket connections it has with other tasks (used in direct message routing), and then closes these connections (step 2'). It is necessary to flush these TCP socket connections to avoid losing any message that may be buffered in these sockets. The details of how the connections



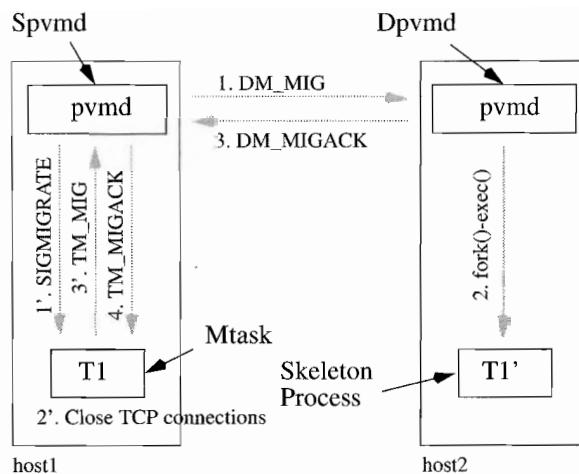


Figure 3. Migration Initialization. Local initializations (1', 2', and 3') and remote initializations (1, 2, and 3) occur in parallel and “sync” at step 4.

are flushed and closed will be discussed in Section 3.3. Note that the tasks with which Mtask had a direct TCP connection continue executing as they normally would. In the event they send a message to Mtask, the message will automatically be forwarded through the pvmds. Message forwarding will be discussed in Section 3.4.2.

Once the TCP connections have been flushed and closed, Mtask sends a *TM\_MIG* CM (TM stands for Task Message) to Spvmd (step 3') to tell the Spvmd that local initialization is complete. Mtask then blocks and waits for a *TM\_MIGACK* CM from the Spvmd.

While the local initialization component is executing in the source machine, remote initialization is proceeding simultaneously on the destination machine. Remote initialization is triggered by a *DM\_MIG* CM (DM stands for Daemon Message) from Spvmd to Dpvmd (step 1, no prime). This CM informs the Dpvmd that a task will be migrating to it. Information about the migrating task such as its tid, executable file name, parent task's tid, etc., is passed along in this CM. The name of the executable file from which the migrating task was started is particularly important since the same executable file is used to start a “skeleton process” (step 2). The executable file is assumed to be accessible from the destination machine. The skeleton process provides the infrastructure to which process state can be transferred and will eventually be executing in the context of Mtask.

State transfer has three requirements: the source of the state, the recipient of the state, and the medium through which the state will be transferred. The first

two components are satisfied by Mtask and the skeleton process respectively. For the transfer medium, a TCP connection, to be established at process state transfer time, is used. The main advantage of using a TCP socket is that TCP provides reliable delivery of sequenced data, simplifying the transfer of data between Mtask and the skeleton process.

For the TCP connection to be established, it is necessary that the skeleton process have a TCP socket to which Mtask can “connect.” Before the Dpvm starts the skeleton process, it creates a socket and binds it to a port address. Following the semantics of `fork()/exec()`, the skeleton process automatically inherits the socket from the Dpvm.

In addition to inheriting the socket, special arguments are also passed to the skeleton process. These special arguments cause the skeleton process to execute “restart code”. Recall that at application start-up, some migration initialization code is first executed prior to executing the application’s code. Part of the migration initialization code is to test whether the process has to execute as a skeleton process, based on the arguments passed to it. If the process was started as a skeleton process, it will wait for a connection on the socket it inherited from the Dpvm. If not, it executes application code.

For Mtask to be able to connect to the socket waited on by the skeleton process, Mtask must know the port address the socket is bound to on the destination machine. This port address is known to the Dpvm. To send the port address to Mtask, the Dpvm sends a *DM\_MIGACK* CM to the Spvm (step 3) containing the port address and an error code. If the error code is zero, then the port address is valid and can be used by Mtask to connect to the skeleton process. A non-zero error code indicates that something went wrong during remote initialization and that migration cannot proceed. Possible error codes are *PvmNoFile* and *PvmOutOfRes*. A *PvmNoFile* error code means that the executable file name of the migrating task was not found on the destination machine. A *PvmOutOfRes* error code means that there wasn’t enough resources on the destination machine to start the skeleton process. This error could be caused by several factors such as inability to create more sockets, inability to `fork()` another process, etc. A non-zero error code causes the Spvm to send the GS an *SM\_MIGACK* CM containing the error code, similar to what it would have done given an invalid tid or hid from an *SM\_MIG* CM. Sending these error codes back to the GS allows the GS to keep track of unsuccessful migrations, giving it an up-to-date view of the state of the PVM system.

The last part of this stage is for the Spvm to send a *TM\_MIGACK* CM to Mtask (step 4). Recall that at the end of the local initialization, Mtask blocks waiting for this message. But before the Spvm can send this CM to Mtask, it must be sure that both the local and remote initializations have completed. Completion

of local and remote initializations is indicated by receipt of both the *TM\_MIG* CM from Mtask and the *DM\_MIGACK* CM from the Dpvmmd.

The *TM\_MIGACK* CM sent to Mtask contains three items: an error code, the IP address of the destination machine, and the port address of the socket to connect to on the destination machine (the one the skeleton process is waiting on). If the error code is zero, then the migration protocol proceeds to the next stage. If the error code is non-zero (for reasons mentioned above), the migration is aborted and Mtask simply returns from the migration handler and continues its execution prior to getting interrupted by the SIGMIGRATE signal.

Conceptually, this stage of the protocol is simple. Unfortunately, the same cannot be said for the actual implementation. In particular, note that while within the migration signal handler, the CMs *TM\_MIG* and *TM\_MIGACK* are sent and received respectively. Sending and receiving these CMs requires the migration signal handler to use routines in the pvmlib. However, the pvmlib is not re-entrant. If the migration signal handler happened to be invoked while Mtask was executing within the pvmlib, the migration signal handler's use of pvmlib routines could corrupt data structures in the pvmlib, leading to unpredictable results.

The obvious solution is to make the pvmlib re-entrant. Making the pvmlib re-entrant, however, would require a complete re-work of the library. The simpler but not so elegant approach of synchronizing migration signal handler invocation with the task's execution in the pvmlib is used. One way of achieving this synchronization is to block the SIGMIGRATE signal whenever task execution enters the pvmlib. The approach was tried and worked as expected. Unfortunately, blocking and unblocking signals require system calls that incur a significant amount of overhead.

The solution used, which gives the same result as that of blocking signals but with much less overhead, is to set an *IN\_LIB* flag whenever task execution enters the pvmlib. When the migration signal handler is invoked, this *IN\_LIB* flag is checked first. If the flag is not set, the migration signal handler can safely use pvmlib routines. If it is set, however, the migration signal handler sets another flag called the *IS\_MIG* flag, indicating that the signal occurred, and returns. When execution of the task leaves the pvmlib, in addition to clearing the *IN\_LIB* flag, the *IS\_MIG* flag is checked. If the *IS\_MIG* flag is set, the task generates a SIGMIGRATE signal to itself. Sending the signal to itself "simulates" the situation when the signal was first received, except that this time, the task is already outside the pvmlib.

A problem arises when execution of the task blocks inside the pvmlib such as when the user code calls a `pvm_recv()` and the desired message has not yet arrived. This situation is undesirable since the migration signal handler will not get

invoked for an indefinite amount of time, preventing the task from migrating. To resolve this problem, all PVM routines that could block for an indefinite amount of time were rewritten so that they blocked outside the pvmlib. That is, modifications were made such that instead of having the routines block deep inside the pvmlib as they would in the original pvmlib, they now block on the “surface” of the pvmlib. This change is sufficient to allow the pvmlib routines used in the migration signal handler to be executed without running into re-entrancy problems.

### 3.2.3. Process State Transfer

Reaching this stage of the migration protocol implies that the skeleton process was successfully started and that Mtask has received the *TM\_MIGACK* CM containing the destination host’s IP address and the port address of the socket the skeleton process is waiting on.

Before the state of Mtask is transferred, Mtask first detaches from the local pvmd (Spvmd in this case) using `pvm_exit()`. This call closes the TCP socket connection Mtask has with its local pvmd. Messages in the pvmlib that have not yet been received by the application remain intact in the task’s data space.

As mentioned above, migration involves capturing the process’s state (text, data, stack, and processor context), transferring it to another host, and reconstructing it. The text of the process can be taken from the executable file from which the process was started. It is for this reason the skeleton process is started from the same executable file from which Mtask was started. Using the same executable file automatically “migrates” the text. The data and stack, however, have to be read directly from Mtask’s virtual memory. As for the processor context, recall that this has already been saved in the stack when the migration signal handler was invoked. By performing the state transfer while within the migration signal handler, coupled with the ability to transfer/restore the stack correctly, the processor context is preserved.

The processor context saved due to the invocation of the migration signal handler contains information regarding where execution should resume in the user’s code. However, if migration is to occur within the signal handler, a second set of processor context information is needed to determine where execution should resume inside the signal handler. Correctly resuming execution inside the migration signal handler is necessary for the signal handler to be able to “return” correctly and restore the process context saved when the signal handler was invoked. For this purpose, a `setjmp()` is called within the migration signal handler just before the actual state transfer. A similar approach is taken in Condor [Litzkow & Solomon 1992].

After calling `setjmp()`, `Mtask` creates a TCP socket and using the IP address and the socket port address from the `TM_MIGACK` CM, establishes a connection with the skeleton process on the destination host. It is through this TCP connection that the data and stack of `Mtask` is transferred.

#### 3.2.4. Restart

After sending all the necessary state information to the skeleton process, `Mtask` terminates. It is at this point where `Mtask` is officially removed from the source host. The skeleton process, after receiving `Mtask`'s state, assimilates it as its own. This assimilation of state is done by placing the received data and stack state in their appropriate place in the skeleton process's virtual address space. A temporary stack is used, again using the signalling facility, while restoring the state of the real stack to avoid corrupting its contents. After restoring all the state information from `Mtask`, a `longjmp()` is done using the buffer saved from the `setjmp()` call in the state transfer stage. This `longjmp()` causes execution to "go back" into the migration signal handler just as it was in `Mtask` at the time the `setjmp()` was called. It is at this point that the skeleton process starts executing in the context of the `Mtask`.

Before the skeleton process could re-participate as part of the application, it first has to re-enroll itself with the local `pvm`d (`Dpvm`d in this case). By re-enrolling to the PVM system, the skeleton process officially becomes an MPVM task, at the same time re-establishing its indirect communications route with the other tasks. As for the TCP connections that were closed prior to the state transfer, note that direct connections are established "on demand" in PVM. That is, only when a message is first sent between two tasks (which have set their routing mode to use direct routing) is the TCP connection established. By closing down the TCP connections in such a way that the tasks involved "think" that there was never a connection, direct connections with the just-migrated task will automatically be re-established, using the protocol provided by PVM, once messages start flowing between them again.

Lastly, though no longer technically part of the restart stage, the `Dpvm`d sends a `SM_MIGACK` CM to the GS containing an error code of zero. This CM informs the GS that the migration was successful and that the migrated task is again up and running.

Figure 4 shows the timeline of the migration protocol. Note that the migration protocol involves only the migrating task, the source `pvm`d, and the destination `pvm`d. Multiple migrations can occur simultaneously without interfering with each other, even if they have overlapping `pvm`ds.

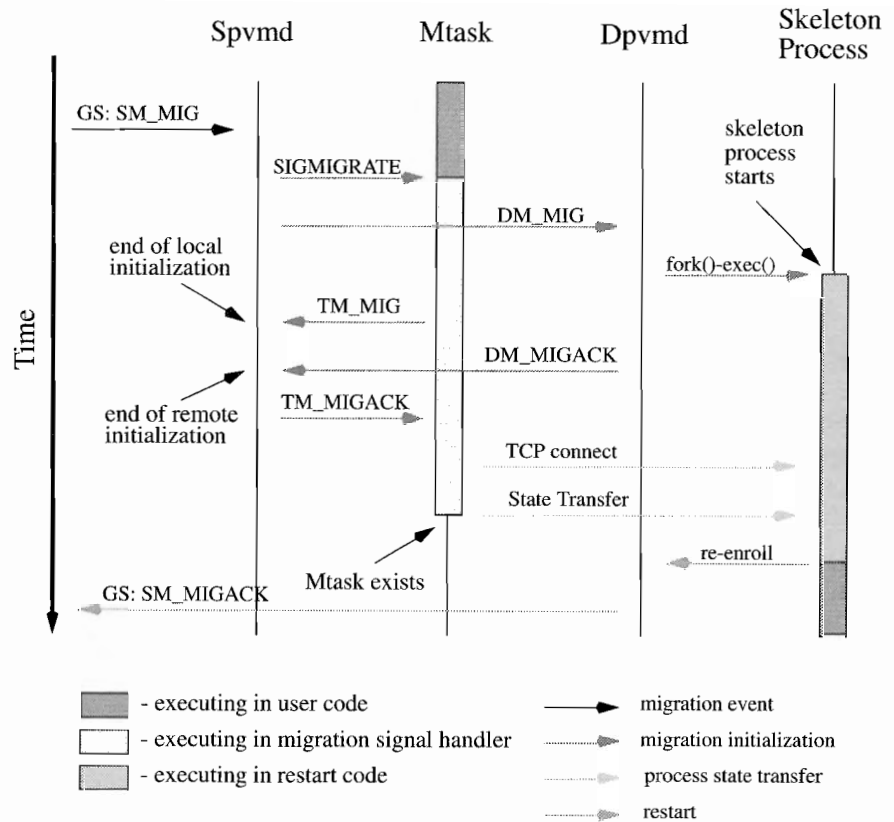


Figure 4. Migration protocol timeline.

### 3.3. Closing TCP Connections

As mentioned in Section 3.2.2, the TCP socket connections Mtask has with other tasks have to be flushed and closed prior to migration. These TCP socket connections are used for direct routing between Mtask and the other tasks. The TCP connections are flushed to avoid loss of any un-received message.

Flushing and closing these TCP socket connections is not as simple as just reading everything that could be read from the socket and then closing them. It is possible that messages are still in-transit and thus not yet available for reading. It is also possible that the peer task (the task at the other end of the connection) is just about to send a message. In either instance, the fact that nothing can be read from the TCP socket connection does not imply that there wouldn't be any in the future.

To ensure that there are no messages in the connection, in-transit or in the future, it is necessary for Mtask to explicitly inform the peer task of its intention of closing the connection and get an acknowledgment from the peer task that it will no longer send messages through that connection. To inform the peer task of the intention of closing the connection, Mtask sends an out-of-band (OOB) data, using the MSG\_OOB flag for the send() system call, through the TCP connection. The OOB data causes a SIGURG signal at the peer task. Using this method of informing the peer task of the connection closure has a number of advantages. First, it enables the peer task to respond to the socket connection closure immediately. Minimizing the time it takes to take down Mtask's TCP socket connection is necessary to minimize the time it takes to migrate Mtask. Second, by testing for "exception conditions" using select(), this method provides enough information for the peer task to know which socket, assuming it also has TCP socket connections with other tasks, is being closed. And lastly, this method involves only Mtask and the peer task which helps minimize the overhead involved in closing the connection.

The exact protocol used is illustrated in Figure 5. TCP socket connections can be pictured as two uni-directional pipes or channels. Another feature of TCP socket connections that the protocol uses is the ability to close only one channel of the connection with the shutdown() system call. Using the shutdown() system call, it is possible to close the TCP socket connection one pipe or channel at a time. The close() system call closes both channels at once.

Initially, Mtask (T1) sends the OOB data to the peer task (T2). It then calls shutdown() to close the channel from T1 to T2, and proceeds to read the socket (i.e., reading the T2 to T1 channel) until it reads an end-of-file (EOF). The OOB data in the mean time causes a SIGURG signal to be generated at the peer task, which in turn invokes a SIGURG signal handler. The SIGURG signal handler at the peer task first determines which socket caused the SIGURG signal using the select() system call, and then reads in all it can from the socket until it reads an EOF. The reading of the socket until an EOF is detected, in effect, flushes any

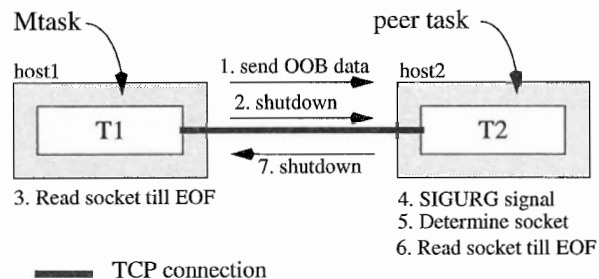


Figure 5. TCP socket connection closure protocol.

un-received messages sent by Mtask. The EOF is guaranteed to be detected due to the `shutdown()` call in step 2. After detecting the EOF, the signal handler issues a `shutdown()` on the T2 to T1 channel. At this point, Mtask is still reading the socket waiting for an EOF on the T2 to T1 channel. Once Mtask reads the EOF, it knows that all the messages, both sent and received through that socket connection, have been flushed.

Three things are worth mentioning about the protocol. First, the protocol works even if T1 and T2 are migrated simultaneously. Both tasks will simply go through steps 1, 2, and 3, with the `shutdown()` on step 2 of one task causing an EOF to be read on step 3 of the other task. Second, note that the SIGURG signal handler, just like the migration signal handler, uses `pvmlib` routines to read the messages from the socket. Thus, the SIGURG signal handler had to be guarded against re-entrancy problems, using the same method used for guarding the migration signal handler. Lastly, the protocol described above assumes only one TCP socket connection is being closed. Typically, either Mtask has no TCP socket connections with other tasks or it has a number of them, all of which have to be closed. Instead of executing the protocol one connection at a time, Mtask can execute steps 1 and 2 for all TCP connections first before going to step 3. This causes steps 4, 5, and 6 on all the peer tasks to be executed in parallel, further minimizing the time it takes to take down all the connections.

A drawback of this method, however, is that it doesn't work for UNIX domain sockets. PVM 3.3.x uses UNIX domain sockets for direct connections between tasks on the same host since it is about 1.5x–2x faster than TCP sockets [Manchek 1994]. UNIX domain sockets only work for tasks on the same host. Unfortunately, UNIX domain sockets have no support for OOB data. As currently implemented, MPVM uses TCP sockets for direct communication even for tasks on the same host.

One last aspect related to TCP connection closure is with regards to routing for messages from the peer tasks to Mtask and the re-establishment of the TCP connection after migration. As mentioned previously, TCP connections between tasks are established on demand. Since the peer task has its routing option set to direct routing (which had to be set in the first place for the just-taken-down TCP connection to have been established), the peer task will try to establish another TCP connection on the next message to Mtask. But since Mtask is migrating, this should not happen. Otherwise, the peer task would only be blocked waiting for an acknowledgment from Mtask. To address this problem, before the SIGURG signal handler returns, a flag inside the `pvmlib` is set to indicate that the peer task should not try to establish a TCP connection with Mtask. Messages for Mtask will then be routed indirectly through the `pvm`s allowing the peer task to continue executing.



Once Mtask has migrated and is running again, it would be desirable for the TCP connections that were taken down before migration to be re-established. Since a flag has been set on the peer tasks, no TCP connection request will come from the peer tasks. The request should come from Mtask. One option would be for Mtask to request a TCP connection from all the tasks it used to have a connection with prior to returning from the migration signal handler. However, this option would pay the price of establishing the connection without knowing if the connection will even be used. This brings connection re-establishment back to the “on demand” philosophy.

To continue supporting “on demand” TCP connection establishment, one possibility is to inform all the peer tasks that they could now establish a TCP connection if they wanted to. This option, however, would require that a message be multicast to all the peer tasks. The solution taken in MPVM currently is to do nothing. That is, a TCP connection will be established with a peer task only if Mtask requests it. This request will be generated on the first message Mtask sends to the peer task after the migration. This implementation, however, implies that if the communication between the peer task and Mtask is always one way from the peer task to Mtask, all the messages will be routed through the pvmds. Both options have advantages and disadvantages. Which one is better is debatable.

### *3.4. Message Delivery on Migration*

An important aspect of the MPVM implementation that has yet to be discussed is how MPVM handles messages for migrating/migrated tasks. That is, how do messages sent to Mtask find their way to the new location of Mtask. To ensure correct delivery of messages in the presence of migration, support for virtual tids, message forwarding, and message sequencing had to be built into MPVM.

Note that the problem of message delivery really only applies to messages using indirect routing. Direct routing is not a problem since by definition it uses a point-to-point connection. Also, recall that at migration time, direct connections are taken down and messages from other tasks to the migrating task are automatically routed indirectly through the pvmds. The next three sections will therefore be presented in the context of indirectly routed messages.

#### *3.4.1. Virtual Tids*

All tasks in PVM are identified by task identifiers or tids. These tids are used to identify the source and destination of messages. Tids are formed using an encoding of a host ID and a task number [Geist et al. 1994]. The host ID or hid

represents the host where the task is executing while the task number identifies a particular task on a particular host.

The combination of the host number and the task number uniquely identifies any task in the entire virtual machine. One advantage of this encoding scheme is that it allows fast routing of messages since the target host of any message can be determined directly from the destination tid. However, recall that the tid of a task is part of the state information maintained on migration. That is, a task with tid T1 will always be referred to as T1 regardless of where it is actually executing. The use of the same tid is necessary to make migrations transparent to the application. Unfortunately, the use of the same tid also implies that there is no longer any guarantee that the host number encoded in the tid is the actual host where the task is executing.

MPVM gets around this problem by virtualizing tids, thus making them location transparent. Virtualizing the tids is done by maintaining a table of tid-to-host mappings. Instead of just relying on the host number encoded on the tid as the search key for the target host, the whole tid is used. Note that the same host number and task number encoding scheme is still used in generating the tids.

Each pvmd in the virtual machine maintains two tid-to-host mapping tables: a *home* map and a *hint* map. The home map on host H, for example, contains a list of mappings for tasks that were originally started on host H, regardless of where they are currently executing. Note that since these tasks were originally started on host H, the host numbers in their tids “point” to host H as their home. The home map on host H is always updated whenever a task whose home is host H migrates. The concept of home/hint maps is not new and have been used in other systems such as V [Theimer et al. 1985] and Amoeba [Mullender et al. 1990].

Consider the example in Figure 6. In step 1, task T1 is started in host H1. This causes a T1 → H1 entry to be added on the home map of H1. At some later time, step 2, T1 migrates to host H2. This migration causes the T1 → H1 home map entry on H1 to be updated to T1 → H2, indicating that T1 is now on H2. The same goes for step 3 when T1 migrates to H3. Notice from the figure that when T1 migrated from H2 to H3, a *DM\_HOMEUPD* CM was sent from H2, where T1 migrated from, to H1, the home of T1 (step 4). This CM informs H1 that task T1 has migrated to H3, causing H1 to update its home map. It was not necessary to have a *DM\_HOMEUPD* CM when T1 first migrated from H1 to H2 since H1 is already the home of T1 and the home map can be updated directly.

With the use of the home map, it is possible to determine the exact whereabouts of any given task. For example, using Figure 6 again, if at some later time task T2 in host H2 sends a message to T1, the message is first routed through the pvmd on H2 (step 5). The pvmd on H2 determines that the destination for the message is T1 and sends a *DM\_HOMEREQ* CM to H1, the home host of T1 (step

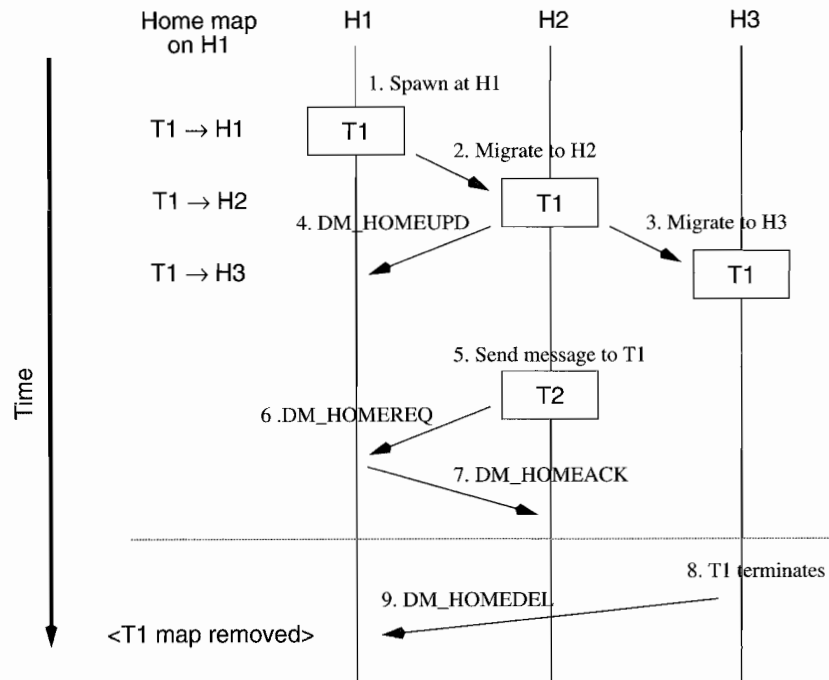


Figure 6. Tid virtualization using home maps. A task's home host is the host where it was started and is identified by the host number encoded in the task's tid.

6). Recall that the home host is determined from the host number encoded in the tid of T1. The pvmd on H1 receives this CM and replies with a *DM\_HOMEACK* CM containing the T1 → H3 mapping (step 7). The pvmd on H2 then knows that the message should be sent to H3.

When T1 terminates, the pvmd on host H3 sends a *DM\_HOMEDEL* CM to the pvmd on H1 indicating that T1 has terminated. This CM causes the home map entry for T1 to be removed. If at some later time a *DM\_HOMEREQ* CM for T1 is received by H1, a *DM\_HOMEACK* CM containing an error code is sent back to the requesting pvmd. This error code would inform the requesting pvmd that T1 no longer exists and the message for T1 is discarded.

While this scheme works nicely, it is terribly inefficient. To improve the performance, a hint map is used. The hint map caches tid-to-host mappings received from previous home map requests. Using a hint map will limit the need for sending *DM\_HOMEREQ* CMs to only those cases when there is no local copy of the mapping. As currently implemented, the hint map is allocated statically. This limits the number of mappings that could be cached. When the hint map gets full, replacement of entries uses the least recently used (LRU) policy. This policy will

throw away the mapping that hasn't been used for the longest time. Examples of such mappings would be those for tasks that have terminated.

To minimize the number of *DM\_HOMEREQ* CMs, home and hint maps in Spvmd and Dpvmd are updated immediately. For example, when task T1 was migrated from H2 to H3, updating the hint map on H2 immediately with the entry T1 → H3 will eliminate the need for a *DM\_HOMEREQ* CM in step 5 of Figure 6. Similarly, hint map entries are also immediately updated for newly spawned tasks.

A problem associated with home maps, however, is that it creates a residual dependency between a task and its home node. That is, the home node has to be functional even if the task is executing elsewhere. If the pvmd on the home node should fail, messages to the task may not get delivered correctly since information about the current location of the task is no longer available.

To resolve this problem, home maps are replicated on other nodes. As currently implemented, the home map list of a node is replicated in two other nodes. The decision of replicating home maps only on a subset of the nodes in the system was for scalability considerations. The ideal case is for a home map to be replicated on all the other nodes.

Each pvmd maintains an ordered list of all the nodes in the virtual machine. The order is determined by the sequence the nodes were added into the system. When a node fails, its failure is detected by the pvmds on the other nodes (e.g., via timeouts) and the node list on each active pvmd is updated to reflect this failure. The back-up nodes of node A are defined to be the next two available nodes following the entry for node A in the node list. Since all the nodes have identical copies of the node list, each pvmd can locally determine the back-up nodes of any given node.

When a pvmd needs to send a *DM\_HOMEREQ* CM, it first checks if the “real” home node is available. If so, it sends the *DM\_HOMEREQ* CM to that host. Otherwise, it will send the *DM\_HOMEREQ* CM to one of the “real” home node's back-up nodes. With the use of back-up nodes, the dependency of a task on its home node is removed. This also implies that a node can be “totally” vacated since pvmds can now be removed, if desired.

A problem with this solution is that if the real home node and its two back-ups should fail, reliable message delivery will again be compromised. To help deal with this problem, when a home node detects that one of its back-up nodes have failed, another node is chosen to replace it. The new back-up node will be the next active node in the nodes list that follows the failed node. In this way, the rule that the back-up nodes are the next two active nodes in the list still applies. By re-electing back-up nodes, a node will always have two back-up nodes (unless the number of nodes is less than three); this minimizes the problem of losing messages to only when the node and its back-ups “simultaneously” fail.

### 3.4.2. Message Forwarding

With the use of home and hint maps, it is possible to determine the exact location of any task at any time. However, in the face of migration, these home and hint maps could be left in an inconsistent state. For example, using Figure 6 again, the home map in H1 won't reflect the  $T1 \rightarrow H3$  mapping until it receives the *DM\_HOMEUPD* CM. If a *DM\_HOMEREQ* CM arrived just before the *DM\_HOMEUPD* CM, the *DM\_HOMEACK* CM reply would contain a  $T1 \rightarrow H2$  mapping which is no longer true. Also, note that nowhere in the migration protocol are the other pvmds (aside from the source pvmd, target pvmd, and pvmd on the home host) in the virtual machine informed of the migrated tasks' new location. Thus, the hint maps on these "uninformed" pvmds could contain old and now invalid tid-to-host mappings. The end result of these invalid home and hint maps is that messages will be sent to the wrong host. In this case, the received message should be forwarded to the correct host and the invalid host/hint maps corrected.

Consider the example in Figure 7. Assuming H1 is the home host of T1, H1 has a home map entry for T1. In step 1, T1, which is currently in H3 (which means task T1 migrated from H1 to H3, possibly through other hosts) migrates to H4. At almost the same time, T2 on H2 sends a message to T1 (step 2). If H2 had an out-of-date hint map, the message would be sent to H3, the previous host of T1. The pvmd on H3 will determine that the destination task T1 is no longer one of its local tasks. At this point, there are two possibilities: either the pvmd on H3 has an idea of where T1 is (it has a home or hint map entry for T1) or it doesn't.

In the case where the pvmd has an idea of where T1 is, H4 in this case, the pvmd on H3 will send the pvmd that sent the message a *DM\_HINTUPD* CM containing a  $T1 \rightarrow H4$  mapping (step 3), and then forward the message to H4 (step 4). The *DM\_HINTUPD* CM will cause the pvmd on H2 to update its hint map so that future messages for T1 will be sent directly to H4. Note that the  $T1 \rightarrow H4$  mapping is not necessarily valid. Such would be the case if T1 migrated again from H4 to some other host. In that case, the message forwarding sequence will simply repeat.

The other case to consider is when the pvmd on H3 has no idea where T1 is currently executing. This case is possible if T1 once migrated to H3, causing a  $T1 \rightarrow H3$  hint map entry to be added on H2, but has since moved to another host and the T1 mapping in H3 has been removed from its hint map due to the LRU policy. Since the pvmd on H3 doesn't know of T1's whereabouts, it sends a *DM\_HOMEREQ* CM to H1, the home of T1 (step 6). It then sends a *DM\_HINTDEL* CM to H2, the source of the message (step 7). Eventually,

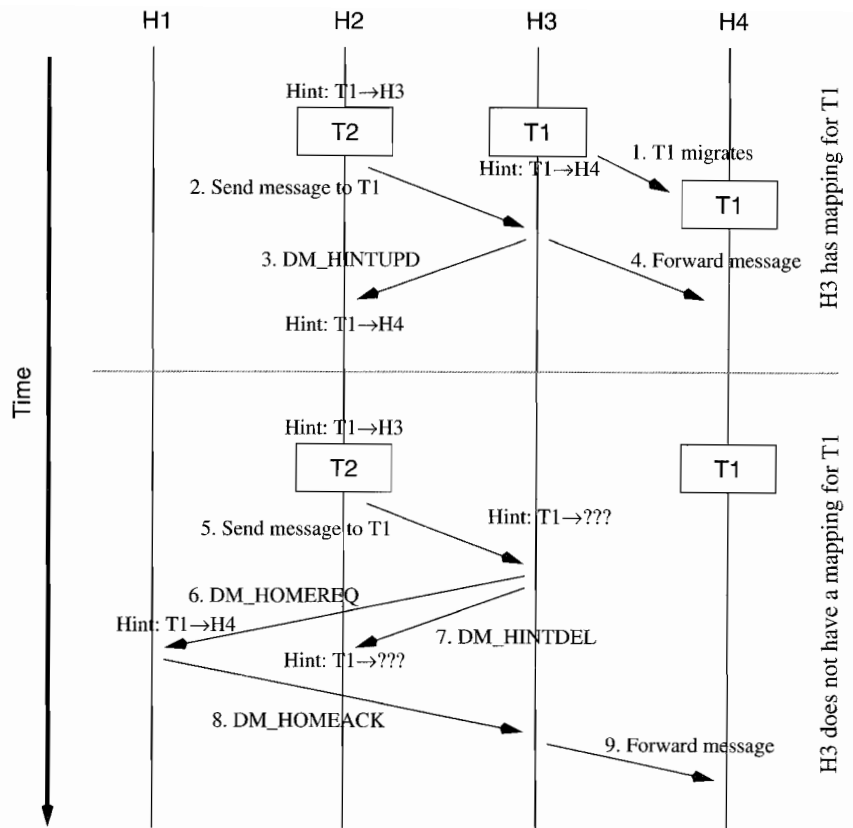


Figure 7. Message forwarding protocol.

the pvmd on H1 will reply with a *DM\_HOMEACK* CM to H3 containing the  $T1 \rightarrow H4$  mapping (step 8). H3 updates its hint map and then forwards the message to H4 (step 9).

The *DM\_HINTDEL* CM sent to H2 in step 6 causes the incorrect  $T1 \rightarrow H3$  hint map entry on H2 to be removed, forcing H2 to request T1's location from H1 on the next message to T1. An alternative implementation is for H3 to wait for the *DM\_HOMEACK* CM from H1 and send the returned mapping to H2 using a *DM\_HINTUPD* CM. This method would update the hint map on H2 eliminating the need for H2 to send a *DM\_HOMEREQ* CM to H1 for future messages to T1. The drawback of this method is that while H3 is waiting for the *DM\_HOMEACK* CM reply to arrive, the pvmd on H2 may be continuously sending messages to H3, all of which have to be forwarded to H4. By sending the *DM\_HINTDEL* CM to H2 immediately, the pvmd on H2 would be forced to get the true location of T1 from H1, allowing the messages to be sent to H4 directly.

### 3.4.3. Message Sequencing

A consequence of message forwarding, however, is that it could break PVM's message ordering semantics. Consider the situation in Figure 8. The situation is similar to the example in Figure 7 above except that the message H2 forwarded to H3 (message A in Figure 8) takes a long time to get to H3. An example of why message A could be delayed is that H2 is on a different network than H1 and H3. Since H1 and H3 are on the same network, a message would travel faster from H1 to H3 than from H2 to H3. The important point here is that the delay, whatever the reason, caused message A to arrive after message B. This behavior is a direct violation of the PVM message passing semantics since message A was sent before message B. It is therefore essential to use some sort of sequencing mechanism to ensure proper ordering of messages.

In standard PVM, the pvmds communicate via UDP sockets for scalability reasons. UDP transport, however, has two basic restrictions. First, a UDP message or datagram can only be UDPMTU (UDP Maximum Transmission Unit) bytes long. The UDPMTU limit is host dependent. This restriction requires messages larger than UDPMTU bytes to be broken up into message fragments or packets. Note that when considering the effective UDPMTU between two hosts, the smaller of the two MTUs is used. For example, the UDPMTU between H1 and H2 is 4096 but only 2048 for H2 and H3. The second restriction

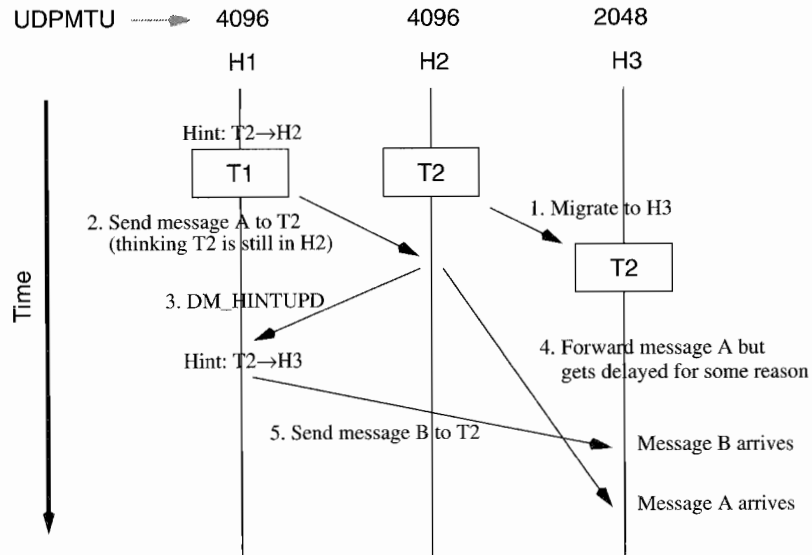


Figure 8. Example where message forwarding breaks message order. In this example, messages A and B are sent from T1 but arrive at T2 in the reverse order.

is that UDP is unreliable. That is, datagram delivery is not guaranteed. These two restrictions requires the pvmds to (1) have the ability to fragment and re-assemble messages and (2) support message fragment re-transmission to guarantee delivery.

Taken in this context, the discussion above regarding message forwarding really applies to message fragments rather than whole messages. Thinking of the example in Figure 8 in these terms (i.e., messages A and B are actually fragments A and B of one message), it is possible that the whole message arrives but with its contents out of order.

To address this problem, each message sent from T1 to T2 is assigned a sequence number based on the number of bytes already sent from task T1 to T2. That is, the sequence number for message  $N + 1$  is calculated as  $Seq_{N+1} = Seq_N + Len_N$  where  $Seq_N$  and  $Seq_{N+1}$  are the sequence numbers of packets  $N$  and  $N + 1$  respectively,  $Len_N$  is the length of packet  $N$  in bytes, and  $Seq_0 = 0$  and  $Len_0 = 0$ .

For example, using the UDPMTU values in Figure 8 and assuming T2 is still in H2, if the first message T1 sends to T2 is 6000 bytes long, it will be broken into two fragments with 4096 and 1904 bytes each with sequence numbers 0 and 4096 respectively. A second 6000 byte message from T1 to T2 will again be broken into two fragments of 4096 and 1904 bytes each but will have sequence numbers 6000 and 10096.

Since each message fragment has a unique sequence number, it is now possible to re-arrange the fragments even if they arrive out of order. Note that even if fragments are further fragmented, correct ordering can still be maintained. Fragmentation of message fragments is possible due to different UDPMTUs between hosts. Using Figure 8 again, the pvmd on H2 had to forward a message fragment from H1 to H3. If the fragment is 4096 bytes long (UDPMTU between H1 and H2 is 4096), it will have to be further fragmented into two 2048 byte fragments since the UDPMTU between H2 and H3 is only 2048. In the re-fragmentation process, new sequence numbers are calculated, using the same equation above, for each of the fragments. For example, if the original 4096 byte fragment had a sequence number of  $S$ , after re-fragmentation, the first fragment will have sequence number  $S$  and the second fragment will have sequence number  $(S + 2048)$ , where 2048 is the length of the first fragment.

Note that message sequence numbers are based on point-to-point messages. That is, the sequence numbers for messages from T1 to T2 are independent of the sequence numbers of messages from T1 to any other task. Since these sequence numbers are based on point-to-point messages, the assignment of sequence numbers and the re-ordering of message fragments based on these sequence numbers is done in the pvmlib. Also note that since the pvmds guarantee delivery of message



fragments through re-transmission, the message re-assembly code in the pvmlib, the code responsible for correctly sequencing fragments into messages, need not worry about lost packets.

Another point to mention is with regards to 0-byte messages (i.e., one produced by a `pvm_initsend()`; `pvm_send()` code sequence). Considering how sequence numbers are calculated, the sequence number of a 0-byte message  $N$  will be the same as the sequence number of message  $N + 1$ . This situation is obviously unacceptable. Fortunately, what the application sees as a 0-byte message is actually a message with some header information and 0 bytes of application data. The message header contains information such as the message's tag and encoding. Since the message headers are counted as part of the message length, there can never be truly 0-byte messages.

While the sequencing mechanism above works for point-to-point messages (i.e., sent by `pvm_send()`), it presents a problem for multicast messages (i.e., sent by `pvm_mcast()`).

Figure 9 illustrates the multicast mechanism in PVM as task T1 sends a multicast message to tasks T2 through T5. In step 1, task T1 first sends a list

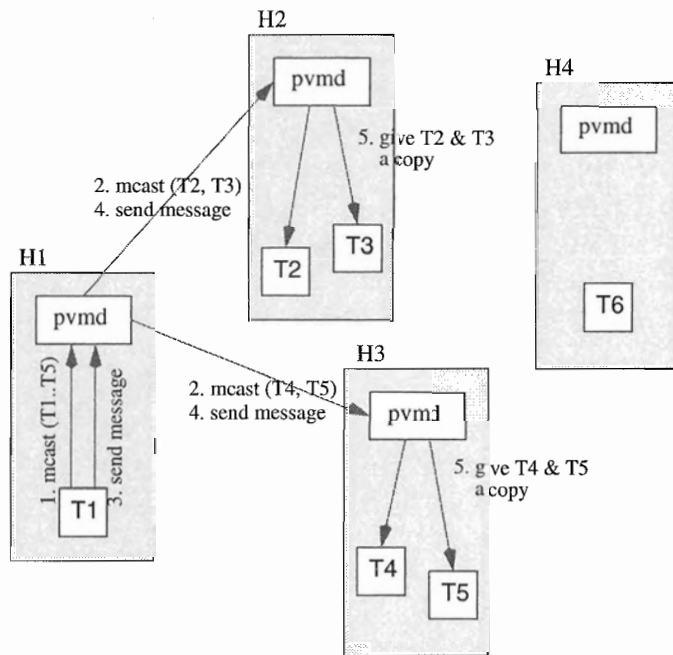


Figure 9. PVM Multicast protocol. This example shows the steps involved in sending a multicast message from T1 to tasks T2 to T5.

of the target tasks (T2 . . . T5) for the multicast message. The pvmd on H1 then determines the hosts where the target tasks are executing, H2 and H3 in this case. The pvmd on H1 then sends a message to H2 and H3 indicating that a multicast message will be sent to all or some of their local tasks (step 2). In the case for H2, the message contains the tids for tasks T2 and T3 indicating that the multicast message will be for tasks T2 and T3 only (assuming there are other tasks on host2). The same goes for H3. Note that no message is sent to H4 since T6 is not a recipient of the multicast message. In step 3, task T1 sends the actual message. A copy of the message is then sent by the pvmd on H1 to H2 and H3 (step 4). The pvmds on H2 and H3, knowing which local tasks the multicast message is meant for from step 2, send each of the target tasks a copy of the message.

The advantage of this implementation is that regardless of the number of target tasks on H2 for example, only one message will be sent from H1 to H2. The pvmd on the target host is responsible for giving each target task a copy of the message. Unfortunately, the fact that only one message is sent out by T1 causes some problem with the sequencing mechanism discussed above. Conceptually, a multicast message to  $N$  tasks is functionally equivalent to a point-to-point message to each of the  $N$  tasks. Since the sequence numbering is based on the number of bytes previously sent from the sending tasks to the target task, each of these  $N$  point-to-point messages will most likely have different sequence numbers. Since only one message is sent from the sending task for a multicast message, there is a problem on how each of the  $N$  target tasks will receive the message with the appropriate sequence number.

This problem is resolved by sending the “would be” sequence number as part of the list of target tasks. Going back to the example in Figure 9, assume that the correct sequence number of the next message for each for the four target tasks T2 to T5 are 200, 300, 400, and 500 respectively. In step 1, instead of sending just the list of target tasks, the list of <tid, sequence number> pairs is sent. The same approach is used in step 2 where instead of just sending a message containing the tids of T2 and T3 to H2, a message containing <T2, 200> and <T3, 300> is sent. When the actual message is sent (steps 3 and 4), the sequence number of the message is reset to 0. That is, if the message is 6000 bytes long and the UDPMTU is 4096, the resulting two message fragments would have sequence numbers 0 and 4096. Finally on step 5, as the pvmd on H2 gives task T2 a copy of each fragment of the message, it adds the sequence number for T2 it got from step 2 to the current sequence number in the fragment. By updating the sequence number, task T2 will receive two message fragments with sequence numbers 200 (200 + 0) and 4296 (200 + 4096), which are the sequence numbers of the next message T2 expects from T1. Task T3 will receive two fragments with sequence numbers 300 and 4396 and so on for tasks T4 and T5. Note that because the sequence number

for the message fragments is reset to 0, this mechanism will work even if the multicast message fragments get further fragmented along the way.

Aside from the modifications described above, more changes had to be made for the multicast mechanism to work under MPVM. Recall that the pvmd on H1 first had to determine the hosts on which the target tasks were executing before it could go to step 2. For this purpose, the home and hint maps are used. It is possible, however, that some of the target tasks don't have an entry on the home or hint maps. One solution is to send a *DM\_HOMEREQ* CM for each task that doesn't have an entry, and only go to step 2 when all the corresponding *DM\_HOMEACK* CMs are received. While this solution would work, it tends to delay messages unnecessarily. For example, if the pvmd on H1 had a mapping for all but T5, tasks T2, T3 and T4 will have to wait until the pvmd on H1 gets the *DM\_HOMEACK* CM for T5.

To avoid unnecessarily delaying the message for all the tasks, instead of sending a *DM\_HOMEREQ* CM for T5, the pvmd on H1 will assume that T5 is in its home host, H3 in this case. If T5 was actually on H3, then everything is fine. If it so happened that T5 was actually on H4, in step 5, the pvmd on H3 will still update the sequence numbers for T5's copy of the message, but will recognize that T5 is not executing locally. This will cause the pvmd on H3 to execute the message forwarding mechanism described previously on the copy of the message for T5. Since H3 is the home node of T5, it knows exactly where the message should be forwarded from its home map. Also, recall that the message forwarding mechanism sends a *DM\_HINTUPD* CM to the source of the message, H1 in this case. This CM will cause the pvmd on H1 to update its hint map so it knows where T5 is located on the next message send or multicast involving T5. This is the same sequence that would happen if the home or hint map on the source of the multicast address had an out-of-date entry for some of the tasks. Such would be the case if the pvmd on H1 had an incorrect T2 → H3 mapping. The only difference is that a *DM\_HOMEREQ* CM might be generated if H3 has no idea where T2 is located.

The use of sequence numbers for multicast messages raises an issue with regards to the message ordering semantics defined by PVM. As mentioned earlier, PVM requires that messages from task A to task B should be received in the same order they were sent. However, this condition actually breaks in PVM when messages are sent using both `pvm_mcast()` and `pvm_send()` with direct routing. Consider the case when a message is sent from task A to task B using `pvm_mcast()` followed by another message sent via `pvm_send()` with direct routing. Since `pvm_mcast()` routes messages through the pvmds while the `pvm_send()` uses a direct TCP connection with task B, it is more than likely that the message sent via `pvm_send()` will get to task B first, in violation of the message ordering semantics. With the use of sequence numbers for both point-to-point and multicast

messages, the message ordering semantics can be preserved. Whether this property gives MPVM some advantage over PVM is hard to say. However, at the very least, mixing `pvm_mcast()` and `pvm_send()` with direct routing will now generate deterministic application behavior.

### 3.5. *Migrating OS State*

OS held state cannot be transferred like the processor state, the process's data or stack. For one thing, since the migration mechanism is implemented at user level, not all OS held state can be captured/reconstructed. An example of process state that cannot be reconstructed is the process ID. Recall that PVM tasks are actually UNIX processes. As such, they have assigned process IDs. Allocation and assignment of process IDs to processes is done entirely by the OS kernel.

Realize that it is only necessary to migrate OS state information that the process can observe directly. For example, OS kernels keep track of the page table entries of processes. But since processes are "usually" not concerned about the specifics of these page table entries, on migration, the OS kernel on the target machine could be left alone in deciding how to allocate pages and page table entries.

The problem in migrating OS state is that the OS state a process observes is valid only in the context of the computing environment at the time the state was observed. Changing the computing environment (e.g., the process migrates from one host to another) would require a mapping of the OS state information as viewed by the process to its equivalent in the new computing environment. This mapping or virtualization of OS state can be achieved to some extent by providing "wrappers" to system calls.

Consider the case of file I/O. To accommodate file I/O migration, the `pvm-lib` supplies its own file I/O routines (e.g., `open()`, `close()`, `dup()`, `read()`, `write()`, etc.) which are wrappers for the actual system calls. The use of wrapper functions in this manner is similar to what is done in Condor [Litzkow & Solomon 1992]. These wrapper functions allow the `pvm-lib` to maintain a list of the files used by a task. This list contains information such as the file's name, file access mode, file descriptor, etc. On migration, but prior to the actual state transfer, for each file in the used files list, the current file pointer offset is taken and then closed. Upon restart, each file in the list is re-opened and the current file pointer is reset to its position prior to migration. The `pvm-lib` also makes sure that each re-opened file is assigned the same file descriptor used before migration.

For file I/O migration to work, as currently implemented, it is necessary that the file be available on the target host. For simplicity, a global file system is assumed to exist (e.g., through NFS). Ways of getting around this restriction are currently being investigated. The current MPVM `pvm-lib` traps only commonly used

file I/O system calls. There is currently no support for `fcntl()` and `ioctl()`, for example.

## 4. Quantitative Evaluation

This section presents performance results for MPVM. The first two experiments were designed to measure the normal case performance (i.e., no migration) of MPVM against PVM 3.3.4 at the micro-benchmark and application level. The third experiment was designed to test migration performance. All experiments that required timing measurements were done on two idle HP series 9000/720 workstations running HP-UX 9.03 connected over an idle 10 Mb/sec Ethernet. Each workstation has a PA-RISC 1.1 processor and 64 MB main memory.

### 4.1. Ping

Table 1 shows the results of running a “ping” experiment using PVM and MPVM. The ping experiment was set up to determine the difference between the message passing times of the two systems. In this experiment, a message is sent from one host to another and back. There is very little computation done. To take the steady state performance for each data size, 50 messages of the appropriate size are first sent back and forth to “warm-up” the system. After which, a timed execution of 1000 ping messages was done.

Table 1. Ping experiment results for PVM and MPVM for direct and indirect communication modes. The numbers represent the average roundtrip time of a message between two hosts.

Data Size	PVM (ms)		MPVM (ms)		% increase of MPVM	
	Indirect	Direct	Indirect	Direct	Indirect	Direct
0	4.788	1.693	4.991	1.990	4%	17%
1024	7.566	4.035	7.760	4.321	2%	7%
2048	9.533	5.432	9.612	5.648	0.8%	4%
4096	17.751	10.059	18.015	10.493	1.4%	4.3%
8192	29.990	18.849	30.024	18.221	0.1%	-3%
16384	53.370	33.694	53.353	33.871	-0.3%	0.5%
32768	102.478	65.211	102.187	65.574	-0.3%	0.5%

As can be seen from Table 1, MPVM and PVM only differ in the order of 10ths of a millisecond. In general, however, MPVM is expected to be slower than PVM for three reasons. First, there is the additional cost of avoiding potential re-entrancy problems. Every time task execution enters/leaves the pvmlib, a flag has to be set/reset. Second, there is the cost of virtualizing tids. This cost only applies for indirectly routed messages. For every message sent out through the daemon, a table lookup has to be done to determine the correct location of the target task. The cost is even greater the first time a lookup is done since a first time lookup would typically result in a *DM\_HOMEREQ* CM mapping request from the home node of the destination task. Lastly, there is the cost of supporting sequence numbers. This sequence numbering cost is linear with the number of fragments of a message.

Notice that as the message size increases, the percentage difference between the performance of PVM and MPVM decreases to the point where they are almost identical. This result indicates that as message size increases, the cost of data transfer increasingly dominates the cost of sending a message.

It should also be mentioned that in the case of direct connections between tasks in the same host, the performance of PVM is better than that of MPVM. This performance difference is due to PVM's use of UNIX domain sockets for direct connections between tasks in the same host as mentioned in section 3.3. MPVM does not use UNIX domain sockets because it doesn't support OOB data which is used for asynchronous closure of direct connections. Ways of getting around this restriction are currently being investigated.

#### 4.2. *Gaussian Elimination*

The Ping experiment above showed the overhead MPVM imposes on message passing performance. Though good for benchmarking, it can hardly qualify as a real-world application since barely any computation was done. To show how MPVM affects the performance of real-world applications, a parallel Gaussian elimination program was run using both PVM and MPVM for different matrix sizes. This experiment only used two machines, each task being responsible for solving half of the matrix. As can be seen from Table 2, the overhead imposed by MPVM is hardly noticeable.

#### 4.3. *Migration Cost*

In this section, the cost of migrating a task is presented. Two measures are defined. The first is the obtrusiveness cost, i.e., the time from when the Spvmd receives an *SM\_MIG* CM to the time the task is removed from the machine (i.e., the

Table 2. Gaussian elimination timing results for PVM and MPVM.

Matrix size	PVM (sec)	MPVM (sec)
80 × 80	0.448	0.453
300 × 300	3.205	3.205
500 × 500	9.311	9.375

task exits). This time represents the minimum time an owner may have to wait before regaining dedicated access to the machine. Note that the machine is not necessarily unusable during this time; it just means that something else is executing other than the owner’s jobs. The second measure is the migration cost. This is the time from when the Spvmd receives the *SM\_MIG* CM to the time the task has restarted on the destination host. The migration cost is essentially the obtrusiveness cost plus the restart stage cost. The first measure approximates the impact of migration on the owner, the other on the job itself. Table 3 and Figure 10 show the obtrusiveness and migration costs for migrating the Gaussian elimination program used in the previous section for various matrix sizes.

Table 3. Obtrusiveness and migration costs for various matrix sizes. The process state size indicates the actual number of bytes transferred at migration time while the TCP transfer time indicates the time spent in sending the appropriate amount of data through a TCP socket connection.

Matrix size	Process state size (bytes)	Obtrusiveness cost (sec)	Migration cost (sec)	TCP transfer time (sec)
0 × 0	97448	0.139	0.327	0.092
80 × 80	109736	0.257	0.361	0.103
300 × 300	277672	0.363	0.590	0.255
500 × 500	597160	0.683	0.871	0.549
1000 × 1000	2100392	1.993	2.205	1.924
2000 × 2000	8109224	7.512	8.324	7.449

As in the previous experiment, this experiment makes use of two machines with one task on each machine. The timing measurements were taken while migrating one task on one machine to the other. The “process state size” indicates the actual number of bytes transferred as measured at migration time. This state size includes static and dynamic data, and the stack. The “TCP transfer time” shows the time spent just transferring the same amount of data over a TCP socket

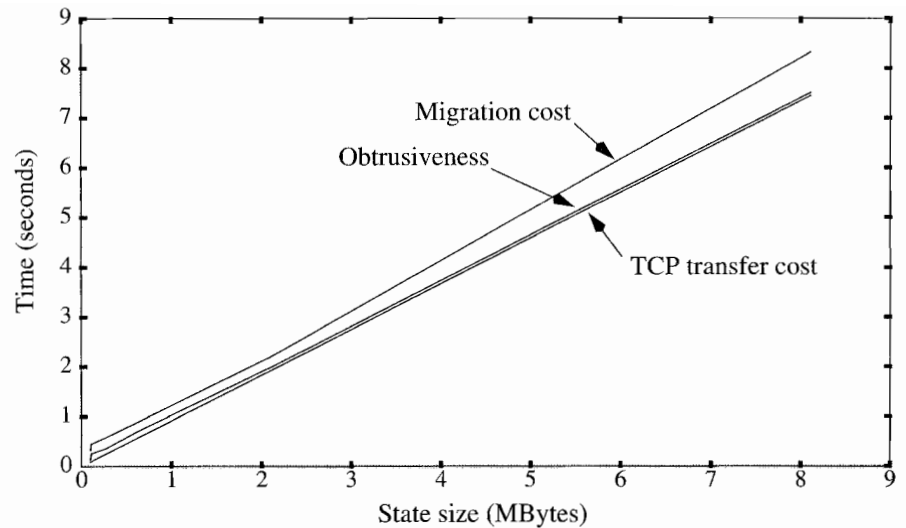


Figure 10. Graph of Table 3. This graph clearly shows the direct relationship between the state size, the TCP transfer cost, and the migration cost.

connection. This measure provides a lower bound for the migration cost. As can be seen from the table, the task's state size is the dominating factor in the obtrusiveness and the migration cost.

While the effect of the migration on the migrating task can be quantified in terms of the migration cost, the effect of migrating one task on the whole application cannot be as easily defined. In the best case, the migration of one task may not affect the performance of the application at all if, for example, the migrated task had the least work to do or was blocked waiting for the other tasks anyway. In the worst case, the entire application could be stalled by as much as the migration cost if, for example, migration occurred just before a global synchronization point (e.g., a barrier), effectively stalling all the tasks in the application until the migrated task resumes execution.

A more important consideration in evaluating the effect of migration on an application is the placement of the migrating task. Migrating a task to a heavily loaded machine will slow down the application not only at migration time but also for the rest of the time it is on the heavily loaded machine. We are currently evaluating the MPVM system using a scheduler [Al-Saqabi et al. 1994] that takes into account processor heterogeneity in terms of architecture and speed. The scheduler recognizes the difference between cooperating tasks (tasks of the same application) and competing tasks (tasks of different applications) and schedules them accordingly to minimize contention. Using the services provided by MPVM, the



scheduler is capable of gang scheduling tasks of multiple applications and allows tasks to double-up on processors. Processor doubling is a technique that allows an application to use fewer resources and yet perform as if all available resources were allocated to it.

## 5. *Related Work*

Process migration implementations can be broadly categorized as either supported at the system-level or at the user-level. In system-level supported implementations, the OS kernel is involved in the migration. Notable examples of such implementation are present in Charlotte [Artsy & Finkel 1989], V [Theimer et al. 1985], Mosix [Barak & Litman 1985; Barak et al. 1993], Sprite [Douglass & Ousterhout 1987; Douglass & Ousterhout 1991], and Mach [Milojicic et al. 1993]. User-level supported process migration implementations, on the other hand, do not require services other than what the OS ordinarily provides through its system call interface. Condor [Litzkow et al. 1988; Litzkow & Solomon 1992; Bricker et al. 1991] is an example of such an implementation. MPVM fits into this category.

Most of the literature on these systems focus on process migration mechanism efficiency. Efficiency is defined in terms of application “freeze time” and state transfer cost. Freeze time is the time during which the process is not executing. State transfer cost on the other hand is the cost of transferring the process’s state. The state transfer cost is often the limiting factor in migration speed [Zayas 1987]. This observation is in accordance with the experimental results shown in the previous section. Two other categorizations of these implementations are in terms of transparency and residual dependency. Transparency refers to how much, or how little, the process is affected by the migration. Of particular interest is how IPC connections are maintained with other processes. Residual dependency refers to how much or how little a migrated process depends on its previous host or hosts.

In Charlotte, the entire virtual address space of a process is transferred at migration time. MPVM uses the same approach. While this approach is simple, it has drawbacks. First, the process is “frozen” for the entire duration of the transfer. Second, the entire virtual address space is transferred even if not all pages may be used by the process. For IPC connections, message senders are informed of the new location of the migrated process. No explicit message flushing is necessary due to kernel provided message caching and retransmission mechanisms. The V kernel addresses the problem of prolonged freeze time by using a technique called “pre-copying.” In essence, while the virtual address space is being transferred, the

process is allowed to continue executing. Once the transfer is complete, the process is stopped, and the memory pages that were touched by the process after the first transfer started are re-copied. The second-stage transfer hopefully is much shorter thus minimizing the freeze time of the process. This technique has been shown to reduce the freeze time significantly though it has to do more work since some pages have to be copied more than once.

While MPVM could benefit from pre-copying, this technique requires access to page table entries—a requirement that cannot be easily satisfied at user-level. As for IPC, the V kernel allows messages to be dropped while the process is frozen. Once the process is un-frozen, the senders are informed of its new location and must re-send the messages. MPVM differs in this regard since MPVM doesn't drop messages. Rather, messages are forwarded to the new destination of a migrating task. If the message arrives at the destination, and the migrating task hasn't restarted yet, the pvmds simply buffer the messages until the migrating task is ready to receive them.

Sprite takes a different approach in transferring process state by making use of a network-wide file system. Since Sprite uses the network file system as a backing store for virtual memory, most of the memory pages used by a process are already saved on the network file system. Hence, at migration time, all Sprite has to do is to flush all the dirty pages of the migrating process and start-up a process on the destination host whose pages are set-up to be demand-paged from the network file system. Mosix uses a slightly different approach for minimizing freeze time by sending all the dirtied pages of the migrating process directly to the target host but sets-up the other pages to be demand paged from the executable file (e.g., the text) or zero-filled. Again, while the techniques used by Sprite and Mosix to minimize freeze time could be used in MPVM, these techniques require access to not only the OS's page table entries, but also the systems swap area (backing store). As for IPC migration transparency, Sprite, Mosix and MPVM use different approaches. Sprite uses the shared file system as the medium for inter-process communication. The file system essentially provides a "well-known" point of communication regardless of where a process is actually executing. In Mosix, transparency is easily achieved since most of the process state information is location independent by design. The obvious disadvantage of MPVM over these two systems is that since MPVM is implemented at user-level, its migration mechanism has to contend with the peculiarities of the OS it is running on as opposed to re-designing the OS to easily accommodate process migration.

Mach uses the concept of Copy-On-Reference (COR) initially used in Accent [Zayas 1987], the precursor of Mach. COR semantics allows a migrating process or task (in Mach parlance) to be started immediately at the target node.

When the new task references a page that hasn't been transferred, a page fault occurs and the page fault handler arranges for the page to be sent from the source host (i.e., demand paging from the source's memory). Compared to MPVM's migration mechanism, this method has the advantage of minimal freeze time and minimal state transfer cost since only those pages actually used by the task are transferred. However, it suffers from residual dependency, due to the fact that resources on the source host cannot be released until either everything has been sent or the task terminates. This residual dependency also makes the process susceptible to failure since if any of the hosts on which the process depends fails, the process could also fail. While MPVM's migration performance may be dwarfed by Mach's use of COR, MPVM doesn't suffer from residual dependencies. To address the transparency problem for IPC connections, a version of Mach that provides in-kernel IPC and DSM called Mach NORMA was used. This method of addressing the IPC transparency issue is very similar to that of MPVM since both systems provide a layer of communication end-point virtualization: the in-kernel IPC in Mach and the message forwarding and task-to-host mappings in MPVM.

On the other side of the implementation domain are those systems implemented at user-level, just like MPVM. The most notable of user-level process migration implementations is Condor. Condor was initially designed for sequential programs. Recently, however, support for PVM applications was added but only as far as scheduling and process suspension/resumption [Pruyne & Livny 1995]. There is currently no support of migration of PVM applications. The main difference between Condor and MPVM is that Condor uses a checkpoint/roll-back mechanism to achieve migration. This design decision was made to minimize obtrusiveness. The Condor system, from time to time, takes a snap shot of the state of the programs it is running. This is done by taking a core dump of the process and merging it with the executable file of the process to produce a checkpoint file. At migration time, the currently running process is immediately terminated. It is later resumed on another host, based on the latest checkpoint file. In addition to being minimally obtrusive, this method has the advantage of fault tolerance in that if something goes wrong (e.g., the system crashes), it is still possible to restart the program from the last checkpoint file. Fault tolerance is something MPVM currently doesn't support. Restarting processes based on roll-backs, however, requires idempotent file operations, a problem MPVM does not suffer from since the state is restarted exactly at the point where it was interrupted. Using roll-backs is particularly troublesome for a parallel/distributed application since it would require synchronous checkpointing of all the tasks or some form of message logging mechanism. MPVM also has the additional advantage of requiring almost no disk I/O. Disk I/O will only occur when transferring pages of the migrating task that

have been paged out. Aside from the speed factor, disk space consumption is also avoided.

Three systems closely related to MPVM are UPVM, Dynamic PVM, and Fail-Safe PVM. UPVM [Konuru et al. 1994], another research effort here at OGI, addresses the problem of the coarse-grained distribution granularity present in MPVM. MPVM migrates tasks at the level of whole processes. UPVM introduces the concept of User Level Processes (ULPs) which are thread-like entities that are independently migratable. Since ULPs are smaller “processing” entities than processes, UPVM has the potential for achieving better load balance. As currently implemented, UPVM has two main restrictions. First, it only runs SPMD programs. Second, since all the ULPs share the address space of a single UNIX process, there is a limit on the number of ULPs the application can have depending on the size of the virtual address space of the process and the memory requirements of each ULP.

Dynamic PVM [Dikken et al. 1994] is an extension to PVM to support process migration, very much like MPVM. The overall designs of Dynamic PVM and MPVM are very similar, such as the use of hint maps and the lazy update of routing information [Dikken 1993]. There are differences, however, in how hints are used. In Dynamic PVM, all hint information is stored and requested from the master pvmd. When a pvmd receives a message for a task that is not executing locally, it replies with a ‘TASK-UNKNOWN’ acknowledgement to the sender of the message. This ‘TASK-UNKNOWN’ acknowledgement causes the first pvmd to request route information from the master pvmd. In this way, Dynamic PVM does not need to forward messages.

More striking differences, however, appear in that Dynamic PVM does not support multicast messages (as defined by PVM), can only migrate one task at a time, and uses Condor-style checkpointing to achieve process migration. In [Vesseur et al. 1995], some performance numbers were reported. For example, migrating a process with 0 and 6 MB of data takes 3.8 and 60.5 seconds on lightly loaded Sun SPARC-2’s. Compare these numbers with those in Table 3.

Finally, Fail-Safe PVM [Leon et al. 1993] is an extension to PVM that implements transparent application checkpointing and restart. The effect of process migration can be accomplished by checkpointing the application and restarting the tasks on different hosts. However, since Fail-Safe PVM wasn’t meant to address process migration, there is a large overhead involved in using Fail-Safe PVM for migration. This large overhead is caused by the need for all tasks to synchronize with each other, flush all outstanding messages, and write their checkpoint files to disk to checkpoint the application.

## 6. Discussion

In this section, a more qualitative discussion on the design and implementation of MPVM is presented.

### 6.1. PVM Source Code Compatibility

Recall that one of the goals of MPVM is to be source code compatible with PVM. To this end, MPVM has maintained the same user-interface, their parameters and semantics, as defined by PVM. The `pvm_sendsig()` routine has to be specially mentioned, however. Since MPVM does not currently support migration of user-installed signal handlers, the use of `pvm_sendsig()` may behave differently for a migrated task that uses signal handlers.

A closely related aspect that affects source code compatibility is the use of a GS. When a GS is used, some of the PVM user-interface calls, `pvm_spawn()` for example, are forwarded to the GS. To maintain full compatibility, the GS should respond to these requests in a PVM compatible way. What the GS does with the requests it receives is outside the control of MPVM.

### 6.2. Portability

Another goal of MPVM is that of portability. This was the motivating factor for choosing a user-level implementation. MPVM was first implemented on HP-PA workstations running HP-UX 9.03. It has since been ported onto SunOS 4.1.3, DEC OSF/1 V1.3, and AIX 3 rel 2.

Although machine dependence of the migration mechanism is unavoidable, the dependence was limited by implementing the migration mechanism using signals, sockets, the `setjmp()/longjmp()` function, etc., all of which are available on most UNIX flavors. Also, no assembly language was used. Everything is written using “C” code.

As long as a process can determine the extents of its data and stack segments at run-time, porting the migration code should not be difficult. Consider the difference between the HP-UX and SunOS versions of MPVM for example. For HP-UX, the following macros are defined

```
#define STACK_TOP ((char *) &stk_var)
#define STACK_BASE ((char *) USERSTACK)
#define DATA_TOP ((char *) sbrk (0))
#define DATA_BASE ((char *) &_data_start)
```

For SunOS, the same macros are defined as

```
#define STACK_TOP ((char *)USRSTACK)
#define STACK_BASE ((char *) &stk_var)
#define DATA_TOP ((char *) sbrk (0))
#define DATA_BASE ((char *) &environ)
```

USRSTACK is a system defined macro which is the absolute address of the beginning of the stack. `Stk_var` is a local variable defined in a function where these macro definitions are used. `&stk_var` thus provides the process an approximate top of stack address which is always more than what is needed to restore, but only as much as the amount of stack space used by a stack frame on a function call. The `__data_start` and `environ` variables define the start of the data space under HP-UX and SunOS respectively. The `__data_start` variable is documented in HP-UX. The `environ` variable on the other hand is not documented but could be determined by using the `nm` UNIX command. And lastly, `sbrk()` is a system call, which when given the parameter 0, returns the address of the top of the heap. Thus, when porting to a new system, only the equivalents of these four definitions need to be determined. In most machines, the usage of `sbrk()` and `&stk_var` should be portable and since USRSTACK is usually defined by the system, this leaves only the value of DATA\_BASE to be determined.

Unfortunately, there are some systems that don't have the USRSTACK macro defined. In this case, the easiest thing to do is to let the process figure out the start-of-stack address at run-time. One way of doing it is to get the address of a local variable declared in the `pvmlib`'s `main()` and "round" that address to the next higher or lower page boundary. Rounding up or down of the address depends on whether the stack grows downward or upward respectively. The resulting address is the start-of-stack address. While this workaround is totally portable, it would fail if the local variable was not allocated on the first stack page. This situation is possible, for example, if enough command line arguments were passed to the process such that it filled up the first page of the stack.

Other potential problems are usually caused by system interface incompatibility. For example, some systems use the `sigvec()` interface to install signal handlers while others use `sigvector()`.

There are special cases however that would require more in-depth investigation. For example, the HP-PA workstations use space registers that contain the addresses of a process's text, data, and stack spaces which are guaranteed to be constant for the lifetime of the process. With migration, however, these addresses are bound to change, and would have to be explicitly updated to the new addresses. Fortunately, the signaling facility in HP-UX (as well as in other OSs) provides a third parameter to the signal handler called the signal context.

This signal context contains the processor state that was saved when the signal was invoked. Using the signal context, the values of space registers can be updated before returning from the signal handler. Another example of this special case is how well the `longjmp()` code interacts with the use of signal handlers that use a temporary stack. Such was the case in the OSF/1 V1.3 port. The `longjmp()` code had safety-checks that detected an error when used with a temporary stack when in fact there is none. Fortunately, the system also provides a lower level `_longjmp()` function that is essentially a `longjmp()` without the error checking.

### 6.3. Transparency

The decision to implement migration at user-level for the sake of portability unfortunately had a negative impact on MPVM's capacity to be truly migration transparent. MPVM can only guarantee transparency for PVM interface calls and some file I/O system calls. Again, there is the assumption that a global-file system is used.

By implementing the migration at user-level, state information managed by the OS kernel such as process IDs and pending signals cannot be automatically preserved on migration. Additional transparency problems appear if the task directly uses UNIX facilities that depend on the location of the task. Examples of such facilities are shared memory, pipes, semaphores, sockets, and shared libraries.

When developing applications of MPVM, special attention has to be given to shared libraries since most compilers/linkers/bundled libraries nowadays are configured to use shared libraries when available. The memory addresses of shared libraries on one machine need not be the same on another. Since shared libraries cache addresses of dynamically loaded modules in the address space of the calling process, migrating the process could cause execution to fail when it tries to call a routine in the shared library on the new machine. For this reason, the developer should explicitly create executable files that are statically linked. This requirement is usually satisfied through some compiler or linker option. Also, recall that the `-Dmain=Main` C compiler flag should also be set.

A possible solution to address this transparency issue in user-level implementations is to provide wrapper functions just like those used for file I/O in Section 3.5 for all system calls. Using these wrapper functions, it would be possible to implement mechanisms that would refer back to the original host for context information about and for the process. Such mechanisms would provide the process with a consistent execution environment regardless of where it is currently executing, similar to what is done in Utopia [Zhou et al. 1992] and Lsbatch [Wang et al. 1993] and to some extent, in Condor [Bricker et al. 1991].

#### 6.4. *Heterogeneity Support*

MPVM supports heterogeneity at the same level as PVM in that processes can be started-up on both homogeneous and heterogeneous architectures. However, migration can only occur within homogeneous machine pools. For example, given ten machines (five Suns and five HPs), a task can be started on each machine. A task on an HP machine however can only migrate to any of the four other HP machines.

The difficulty in supporting heterogeneous migration is that process state on heterogeneous machines is represented differently. Heterogeneity can come in the form of different processors and instruction sets, different OSs, different memory management units, etc. Translation of a process's state as captured on one machine to one of a different architecture is not easy, though there is some work being done that addresses this problem [Theimer & Hayes 1992].

#### 6.5. *Scalability*

Recall that the migration protocol generally only involves the migrating task, the source pvmd, the destination pvmd, and the home pvmd of the migrating task. This approach implies that regardless of the number of tasks on the system, the operations required to migrate a task remain the same. The involvement of other tasks would only depend on whether they have TCP connections with the migrating task that have to be closed. Other than that, all other tasks will continue executing as they normally would and will only get affected by migration if they require a message from the migrating task (i.e., application level synchronization). For this reason, the migration protocol is scalable with respect to the number of tasks.

Also, note that the migration of one task is totally independent of the migration of another. This "independence" property of the migration protocol allows multiple simultaneous migrations to occur. Thus, it doesn't matter whether machines are reclaimed by their owners one at a time, all at the same time, or are reclaimed/released at arbitrary times. In this regard, the migration protocol is scalable with respect to the dynamics of the shared network of workstations.

A factor that negatively affects the migration protocol's scalability, however, is the assumption of the existence of a global file system. MPVM currently relies on a global file system in two ways. First, to avoid moving the migrating process's text, it assumes the executable file of the migrating process is available on the destination machine. Second, the current support for migration transparent file I/O assumes that files available on the source machine are also available on the destination machine. Both of these assumptions will only be always true under



a global file system. While such global file systems already exist, the reality is that such file systems are not yet commonplace, though it is certainly possible to “simulate” one via NFS, for example.

## 6.6. Performance

As has been concluded by other studies and also from the task migration cost measurements in Section 4, the performance of the migration mechanism is largely dependent on the cost of transferring the process’s virtual address space. The first implementation of MPVM used a checkpoint file style of process migration. This method requires five disk accesses (each proportional to the process state size): write the core file, read the core file, read the executable file, write the checkpoint file, and then a final read of the checkpoint file when it is restarted on the destination machine. Disk I/O is particularly troublesome if the local file system is actually an exported file system (e.g., via NFS) since it will not only generate lots of disk I/O on the exporting machine but also a lot of network traffic. By moving to direct state transfer through a TCP connection, the migration speed was increased approximately 10x for processes that use lots of memory.

The current implementation of MPVM is similar to that of Charlotte. That is, the entire virtual address space (data and stack at least) is transferred at migration time. As mentioned previously, this has two drawbacks. The first is prolonged freeze time and the second is possible waste of work by transferring all the pages in the virtual address space even though not all may be used.

Unfortunately, current OSs don’t leave much of a choice as far as user-level implementations are concerned. The solutions presented by systems such as V and Mach rely on virtual memory functions such as trapping page faults, checking for dirty pages, etc. These functions, however, are not generally available at user-level. Though there is work being done to provide user-level virtual memory management [Appel & Li 1991; Harty & Cheriton 1992; Sechrest & Park 1991], until such functionality becomes widely available, portable user-level process migration implementations cannot make use of methods available to system-level implementations.

## 7. Conclusion

MPVM is an extension to PVM that provides for transparent process migration. Such a facility allows tasks to be scheduled on a machine and then later moved to another if so desired. This ability to move tasks makes it possible to use idle

cycles on available machines and at the same time respect ownership of those machines. As is, existing PVM applications can be used under MPVM with little modification. Migration is transparent to the application developer as far as the PVM interface is concerned. File I/O migration is also supported to some extent. Versions of MPVM currently exist for HP-UX 9.03, SunOS 4.1.3, DEC OSF/1 V1.3, and AIX 3 rel 2. Micro-benchmarks show that message-passing in MPVM is just slightly slower than that of PVM. However, tests with real-world applications such as the Gaussian elimination program where some amount of computation being done show that this difference in latency is barely noticeable.

To ensure that task migration doesn't affect the correctness of the application, a strict migration protocol is used. The protocol ensures that messages are not lost and are received in the correct order. The design of the protocol is scalable such that the migration of a task is not affected by the number of tasks in the system and multiple simultaneous migrations can occur. The current limitation of the protocol is the assumption of the use of a global file system.

Measurements of migration costs show that the dominant factor in the migration time is the transfer of the process's virtual address space through the network. This bottleneck has been addressed by system-level process migration implementations. Unfortunately for user-level implementations, unless the OS provides user-level memory management functionality, it would seem that nothing else can be done to improve migration performance.

Though the migration mechanism requires processes to be frozen for some time, the important thing to realize is that this very same mechanism allows PVM applications access to machines they couldn't have used otherwise. It is now possible to have long-running applications execute on a more powerful virtual machine owned by someone else without worrying about getting in the way of the owner. Also, machine owners will likely allow use of their machines knowing they will regain dedicated access whenever they want it. Thus, despite the cost of migration, the ability to migrate could lead to large gains in overall resource availability and performance.

We are currently using MPVM within the research group to get some practical experience regarding its usefulness and performance. Combined with a scheduler capable of dynamic task scheduling and gang scheduling, we are currently investigating the real effect of task migration on both the PVM application and the workstation owners.

As for future work, a lot of things still have to be done to improve migration transparency: non-reliance on a global file system, support for migrating user-installed signal handlers, use of UNIX domain sockets for direct communication between tasks on the same host, etc. Support for migrating applications using X-windows will also be investigated. Another aspect being considered is the support

for fault-tolerance with the use of checkpointing. Integration with existing utilities such as batch schedulers (Condor and DQS [Green & Snyder 1993]), tools (Ptools [Gropp & Lusk 1994]), profilers and debuggers (XPVM [Kohl & Geist 1994]), etc., is also being considered. All this work is targeted for the next generation PVM system.

## References

1. K. Al-Saqabi, S. W. Otto, and J. Walpole, Gang scheduling in heterogenous distributed systems, Technical Report, Dept. of Computer Science and Engineering, Oregon Graduate Institute of Science and Technology, 1994.
2. A. W. Appel and K. Li, Virtual memory primitives for user programs, *Proceedings of the 4th International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 96–107, Santa Clara, CA, 1991.
3. Y. Artsy and R. Finkel, Designing a process migration facility—the Charlotte experience, *Computer*, 22(9):47–56, 1989.
4. A. Barak, S. Guday, and R. G. Wheeler, *The MOSIX Distributed Operating System—Load Balancing for Unix*, Lecture Notes in Computer Science. Springer-Verlag, 1993.
5. A. Barak and A. Litman, MOS—A multicomputer distributed operating system, *Software—Practice & Experience*, 15(8):725–737, 1985.
6. A. L. Beguelin, J. J. Dongarra, A. Geist, R. J. Manchek, S. W. Otto, and J. Walpole, PVM: Experiences, current status and future direction, *Supercomputing '93 Proceedings*, pages 765–6, 1993.
7. A. L. Beguelin, J. J. Dongarra, A. Geist, R. J. Manchek, and V. S. Sunderam, Heterogeneous network computing, *Sixth SIAM Conference on Parallel Processing*, SIAM, 1993.
8. A. Bricker, M. Litzkow, and M. Livny, Condor technical summary, Technical Report, University of Wisconsin at Madison, 1991.
9. L. Dikken, DynamicPVM: Task migration in PVM, Technical Report, TR-ICS/155.1, Shell Research, 1993.
10. L. Dikken, F. van der Linden, J. Vasseur, and P. Sloot, *Dynamic PVM—Dynamic Load Balancing on Parallel Systems*, volume Proceedings Volume II: Networking and Tools, pages 273–277, Springer-Verlag, Munich, Germany, 1994.
11. J. J. Dongarra, A. Geist, R. J. Manchek, and V. S. Sunderam, Integrated PVM framework supports heterogenous network computing, *Computers in Physics*, 1993.
12. F. Douglis and J. Ousterhout, Process migration in the Sprite operating system, *Proceedings of the 7th IEEE International Conference on Distributed Computing Systems*, pages 18–25, Berlin, West Germany, 1987.
13. F. Douglis and J. Ousterhout, Transparent process migration: Design alternatives and the Sprite implementation, *Software—Practice & Experience*, 21(8):757–785, 1991.
14. A. Geist, A. Beguelin, J. Dongarra, W. Jiang, R. Manchek, and V. Sunderam, PVM 3 user's guide and reference manual, Technical Report ORNL/TM-12187, Oak Ridge National Laboratory, 1994.
15. T. Green and J. Snyder, DQS, a distributed queuing system, Technical Report, Supercomputer Computations Research Institute, Florida State University, 1993.

16. W. Gropp and E. Lusk, Scalable Unix tools on parallel processors, *1994 Scalable High-Performance Computing Conference*, pages 56–62. IEEE Computer Society Press, 1994.
17. K. Harty and D.R. Cheriton, Application-controlled physical memory using external page-cache management, *Proceedings of the 5th International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 187–197, Boston, Massachusetts, 1992.
18. J. A. Kohl and G. A. Geist, XPVM: A graphical console and monitor for PVM, *2nd PVM User's Group Meeting*, Oak Ridge, TN, 1994.
19. R. Konuru, J. Casas, S. Otto, R. Prouty, and J. Walpole, A user-level process package for PVM, *1994 Scalable High-Performance Computing Conference*, pages 48–55. IEEE Computer Society Press, 1994.
20. J. Leon, A. Fisher, and P. Steenkiste, Fail-safe PVM: A portable package for distributed programming with transparent recovery, Technical Report CMU-CS-93-124, Carnegie Mellon University, 1993.
21. M. Litzkow and M. Solomon, Supporting checkpointing and process migration outside the Unix kernel, *Usenix Winter Conference Proceedings*, pages 283–290, San Francisco, CA, 1992.
22. M. J. Litzkow, M. Livny, and M. W. Mutka, Condor—A hunter of idle workstations, *Proceedings of the 8th IEEE International Conference on Distributed Computing Systems*, pages 104–111, 1988.
23. R. Manchek, PVM version 3.3.0 release-notes, <ftp://netlib2.cs.utk.edu/pvm3>, 1994.
24. D. S. Milojevic, W. Zint, A. Dangel, and P. Giese, Task migration on the top of the Mach microkernel, *MACH III Symposium Proceedings*, pages 273–289, Sante Fe, New Mexico, 1993.
25. S. J. Mullender, G. van Rossum, R. van Renesse, and H. van Staveren, Amoeba—a distributed operating system for the 1990s, *IEEE Computer*, 23(5):44–53, 1990.
26. J. Pruyne and M. Livny, Providing resource management services to parallel applications, *Proceeding of the 2nd workshop on Environments and Tools for Parallel Scientific Computing*, pages 152–161, 1995.
27. S. Sechrest and Y. Park, User-level physical memory management for Mach, *Mach Symposium Proceedings*, pages 189–199, Monterey, California, 1991.
28. M. M. Theimer and B. Hayes, Heterogeneous process migration by recompilation, Technical Report CSL-92-3, Xerox Palo Alto Research Center, CA, 1992.
29. M. M. Theimer, K. A. Lantz, and D. R. Cheriton, Preemptable remote execution facilities for the V-System, *Proceedings of the 10th ACM Symposium on Operating Systems Principles*, pages 2–12, Orcas Islands, Washington, 1985.
30. J. J. J. Vasseur, R. N. Heederick, B. J. Overeinder, and P. M. A. Sloot, Experiments in dynamic load balancing for parallel cluster computing, *ASCI 1995 Conference Proceedings*, accepted for publication, 1995.
31. J. Wang, S. Zhou, K. Ahmed, and W. Long, Lsbatch: A distributed load sharing batch system, Technical Report CSRI-286, Computer Systems Research Institute, University of Toronto, Canada, 1993.

32. E. R. Zayas, Attacking the process migration bottleneck, *Proceedings of the 11th ACM Symposium on Operating Systems Principles*, pages 13–24, Austin, Texas, 1987.
33. S. Zhou, J. Wang, X. Zheng, and P. Delisle, Utopia: A load sharing facility for large, heterogeneous distributed computer systems, Technical Report CSRI-257, Computer Systems Research Institute, University of Toronto, 1992.