

Migration and Rollback Transparency for Arbitrary Distributed Applications in Workstation Clusters^{*}

Stefan Petri^{1**}, Matthias Bolz², Horst Langendörfer²

¹ Institute for Computer Engineering, Medical University at Lübeck
petri@iti.mu-luebeck.de

² Institute for Operating Systems and Computer Networks,
Technical University Braunschweig

Abstract. Programmers and users of compute intensive scientific applications often do not want to (or even cannot) code load balancing and fault tolerance into their programs.

The Φ BEAM system [18] uses a global virtual name space to provide migration and rollback transparency in user space for distributed groups of processes on workstations. Applications always use the same virtual names for the operating system objects, independent of their current real location. The system calls are interposed and their parameters translated between the name spaces. Unlike other migration mechanisms, Φ BEAM does not require the applications to be written for a specific programming model or communication library.

The first approach to execute applications in the virtual name space was to link the programs with a modified system library. Now, In this paper we describe design and implementation of a separate system call interposition process [3] that accesses the application via the debugging interface. The main advantage of this approach is that it can handle even unmodified (e. g. commercially bought) application programs. We compare measured performance figures with previous similar approaches [15, 20] and the modified system library.

1 Motivation and Introduction

Networks of Workstations (NOWs) become increasingly attractive as platforms for parallel compute intensive applications, because their price/performance ratio is significantly better than that of massively parallel systems (MPPs) [10].

However, in contrast to MPPs most NOWs operate in multi user mode, the nodes are shared between applications, and interactive users may not be disturbed by resource hungry background computations. These constraints make it desirable to use a dynamic load balancing facility that moves work from overloaded to idle nodes, or evacuates machines that are claimed for other purposes, e. g. interactive users. Additionally, the high probability of machine failures [14, 6] necessitates fault tolerance measures for long running applications.

Programmers and users of compute intensive scientific applications are mainly interested in getting their problem solved. They expect load balancing and fault tolerance as services of the underlying operating or run time system and do not want to care

* extended CD-ROM version

** At the time of writing funded by DFG contract SFB 342 at Institute for Computer Science, Technical University Munich

about them in their application code. Unfortunately, off-the-shelf workstation platforms provide these services only minimally, if at all.

In the next section, we give a brief overview of our Φ BEAM load balancing system for distributed applications on clusters of workstations. Because the basic problems of “freezing” the application’s state and reviving it afterwards is the same in both cases, checkpointing or migration, Φ BEAM, like other systems (e. g. [13, 22]), handles them both. After that we focus on our approach to apply Φ BEAM transparently to unmodified binary application programs, and show some performance figures. We conclude after a short comparison to related work.

For a broader discussion of concepts and more technical details beyond this paper we refer to [18, 20, 17].

2 Overview of Φ BEAM

The goals of this project are to provide application transparent process migration and checkpointing / rollback for distributed applications on clusters of workstations, running in user space on unmodified Unix systems.

The state of a distributed computation consists of the states of its processes and the states of the communication links between them. We distinguish between a process’ *internal* state (address space and register contents), and its *external* state, the process’ relation to the world outside its address space, the allocated resources, related processes and the communication peers. A process can manipulate its external state only through services of the system kernel. These services operate on objects like files, processes, communication endpoints, etc. They are invoked through system calls, whose arguments name the objects to operate on (file names, process numbers, transport addresses, etc.).

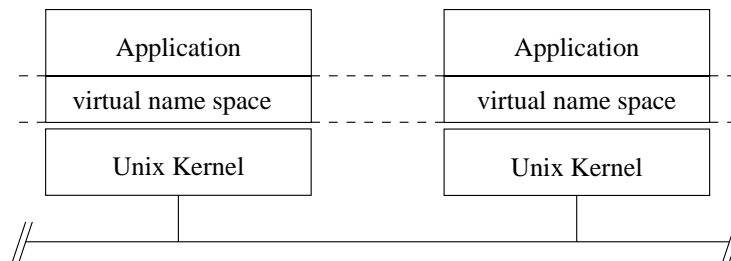


Fig. 1. The virtual name space between kernel and application.

These names depend on the current location and time of the process execution. If a process is moved from one node to another, or an application is restarted from a checkpoint, the names for the objects will change. However, to provide transparency, the application must be able to work with the same names regardless of being moved in space or time. Therefore, Φ BEAM introduces a system wide virtual name space for process IDs, transport addresses and file names. A virtual process table and virtual address tables are maintained for mapping between the virtual and current real names. The applications’ system calls are interposed, the parameters are translated from the virtual into the underlying current real name space, and then the real system service is executed. In the

same manner the return values are translated back from the real into the virtual name space. Through the system call interposition also the changes to the external state are tracked.

An early design decision was to not do any system kernel modifications [18, 9]. Basically, there are two possibilities to perform the system call interposition outside the kernel. Our first approach was a modified system call library. The system call functions are replaced by or wrapped into our own versions [4, 13, 9]. It requires that the application can be linked with the modified library.

The other possibility is to control the application processes by a separate controller process via the debugging interface. The main advantage of this approach is that it can work with applications that are available in binary form only, already completely linked, e. g. commercially bought programs.

In the \wp B/EAM system we have separated the system call interposition component from the name space administration (\wp B/EAM demon `pbeamd` in fig. 2), thus making it easy to switch between different versions of both. To explore scalability, we have also implemented a centralized and a distributed version of the name space administration [19].

The system call interposition component has three main tasks: (i) During normal operation, it does the name space translations, as described above. (ii) When doing a checkpoint or migration, it has to capture the internal and external state, save it, or transfer it for a migration. (iii) Read a saved state from disk or from a network connection and reconstruct a running process from this state information.

3 Operation of the System Call Interposition Process (SCIP)

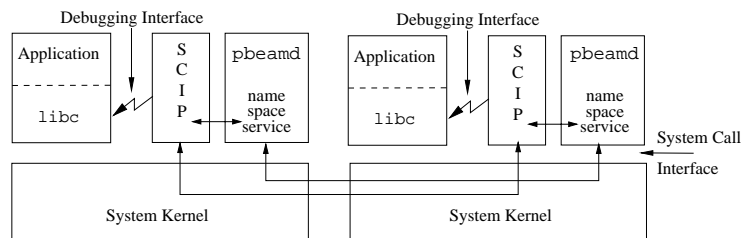


Fig. 2. The SCIP in \wp B/EAM.

The debugging interface offers services to access the internal state of another process and to control its execution. We exploit this to transparently let applications run in our virtual name space. The original Unix debugging interface `ptrace()` only allows the inferior process to be halted at breakpoints, continued, and the memory contents be read or written word-wise. These services are sufficient to save the internal state for checkpointing, although they are inefficient. They are insufficient for performing the described name space translations. However, in most modern Unix versions the debugging interface was extended by services to access entire ranges of the inferior's memory, and to let it stop at entry and exit of system calls. Also they offer extensions that can help to overcome the performance bottlenecks that have been reported in previous work

for similar approaches [15]. These extensions range from enriched `ptrace()` variants to different kinds of process file systems and combinations of both [23, 7].

In the following subsections, after a short discussion of the eligible debugging interfaces, we take a closer look at the problems raised by the above mentioned three tasks of the SCIP, and their possible solutions.

3.1 Controlling Another Process

Since the application program – at least in the beginning – contains neither code for exchanging information with its control process nor for being controlled, SCIP depends on the services the system kernel provides for debugging. The only debugging facility found on all Unix systems is the `ptrace()` system call. But the standard Unix `ptrace()` facility offers only minimal support for breakpoint debugging. So using the `ptrace()` interface for SCIP is only possible by using machine dependent extensions to this call. This makes the implementation more difficult and less portable than the modified system library. Besides that it imposes bad performance: each `ptrace()` request causes a context switch to the controlled process, while traditionally the amount of data being transferred by one `ptrace()` call is only one word.

Another standardized Unix debugging interface is the `/proc` file system defined in System V Release 4 [7]. This virtual file system provides access to another process' address space by means of standard `read()` and `write()` operations and a variety of `ioctl()` requests to control its execution. It promises better performance, more sophisticated control and a good chance of being portable, but a complete implementation of a `/proc` file system is only rarely (if at all) found in the Unix systems in use today.

In the design of the SCIP, we have considered the specifics of SunOS, Solaris, Linux and Irix. Currently, our main implementation platform is a cluster of workstations running SunOS 4.1.x, which does not have a `/proc` file system but several `ptrace()` extensions. Work on Solaris 2.x has just begun.

3.2 System Call Interposition

In our terms, the system call interposition consists of several phases:

1. *Executing the process until the next system call entry.*
2. *Notification of a system call entry.* It must be stopped *before the system call is performed*. The system call arguments including the system call number must be read.
3. *Changing the arguments*, i. e. the virtual names among the arguments must be translated into the real name space.
4. *Performing the system call.*
5. *Notification of the system call exit* and its return values, i. e. the process has just exited kernel mode.
6. *Changing the return values* again from real to virtual names.
7. *Restarting the last system call* if necessary.

As noted before, steps 2, 4 and 5 are supported by most modern Unix flavors. SunOS offers a special `ptrace()` request indicating each system call entry and exit. The `/proc` file system even allows for naming a set of system calls to be observed, leaving uninteresting calls untouched. The popular debugging tools `strace` and `truss` are based on this feature. The important difference between SCIP and a tool that just wants to show which system calls are executed with which arguments and results, is that the former also wants to perform the system call with changed arguments, or even not to perform it at all. This is a task the debugging interfaces obviously were not designed for.

On SunOS, when SCIP gets notified about a system call entry, the arguments cannot be changed. Also there is no way of forcing the process to leave kernel mode without performing the system call. The implementation via the `/proc` file system does not encounter this problem.

On the other hand, changing the results after the process has left kernel mode simply works as one would assume. Also, if an argument in fact is a pointer to a memory area in user address space holding the actual data, then the pointer cannot be changed, but the memory contents can. This is true, e. g., for socket addresses.

For these reasons we distinguish between several classes of system calls. In the simplest case the arguments can (or do not have to) be changed in user space, the traced process may safely perform the system call and the return value is adjusted. Examples in this class are `getpid()` or `gettimeofday()`.

Other calls, like `sigblock()`, cause effects that can and must be changed or reverted afterwards, in most cases by placing some code into the inferior's address space and forcing it to execute this code, thus patching the application at run time.

Things get even more difficult with non-idempotent calls causing unwanted effects that cannot be reverted, e. g. most I/O operations. Here we must be a bit inventive. In case of I/O operations, successful execution can be prevented by immediately closing any file (or socket) descriptor each time one is opened by the application. Subsequent operations on this descriptor will fail, giving SCIP the opportunity to perform the actions with translated arguments. Because doing I/O operations on the SCIP side involves transferring relative large amounts of data between SCIP and the inferior, we will discuss several optimizations along with the performance figures in section 4.

If the application forks a child process, also a new SCIP is spawned that attaches itself to the application's child.

3.3 Freezing the Application

After the `ΦBEAM` demon has requested SCIP to initiate a checkpoint or migration, the traced process must be stopped. If a system call is in progress, that must be finished first, possibly interrupting a blocking call like `sigpause()`. Such an unintended interrupted system call must be repeated when restarting the checkpoint.

Next all the external state information must be read that was not already recorded by the interposition of system calls, e. g. the contents of the interval timers, the set of pending signals and installed signal handlers. Only some of this could be read from the process' user page using a `ptrace()` extension, so SCIP again uses a relocatable compiled function that the traced process is forced to execute.

At last the process' address space must be saved, i. e. written into a file descriptor. One possibility is to copy the whole address space piece by piece into SCIP's address space, then write it. Even if there is a `ptrace()` extension as in SunOS for transferring a large piece of memory with only one call, there is a remarkable overhead, which should be avoided – especially if we consider that migration is most probably initiated when the current machine is overloaded.

A more complex, but faster alternative is to use the Unix file descriptor passing feature. SCIP passes the access rights to its checkpoint/migration file descriptor to a function again patched into the inferior process. This can then efficiently write the whole address space down from “inside” without any detouring.

A significant part of information about the external state is contained in SCIP's own data structures. Therefore, SCIP in fact does a checkpoint of, respectively migrates, itself using the checkpoint mechanism provided by our modified system library. In the middle of freezing itself, it freezes the application state and hides the second process image in its own checkpoint.

To save a checkpoint, a new executable file is built from the text and data parts of the process, while the stack contents are written into a separate file. The files are written either locally or via TCP to a server process. In case of migration, the server process does not write the data to files, but overlays its own address space with the executable image sent on TCP, and then passes control to the restart mechanism.

3.4 Restarting the Application

In both cases, whether restarting from a checkpoint executable or after migration, the first step is to revive the saved SCIP image on the new host, which then in turn restarts the application process from the data contained in its address space image. To restore the application's stack contents, a function containing a call to the `read()` kernel service is patched into the running process.

Thus when restoring the application's process image no (slow) memory transfer via the debugging interface is necessary. This is rather important because the `ptrace()` extension in the used SunOS version that should write larger pieces of memory from the controlling to the controlled process does not work correctly, and the alternative of poking the stack image with `ptrace()` word by word into place has shown to be unacceptably slow (> 10 min / 8 MB on Sparc IPC).

After the application is restarted, the SCIP restores the I/O state, that is open files and network connections. Once more the signal status and the interval timers are restored using a relocatable function patched into the application. Signals that had been sent to the process before checkpointing but had not been delivered because they were blocked, are resent to the process with `kill()`.

After SCIP has restored the saved set of registers via the debugging interface, it continues with controlling the application exactly where it had been interrupted.

Note that, while all that function code, that we patch into the controlled process, is of course machine dependent, the *functionality* of them is not. In order to construct them merely a compiler that can produce relocatable code (e. g. GCC) and several assembler instructions encapsulated in preprocessor macros are needed. In fact, the use of code patching instead of accessing kernel data structures decreases porting effort.

The architecture, operation principle and most of the code of SCIP can be kept when porting.

4 Performance Figures

In this section we summarize performance for interposing selected system calls and the costs of migration/checkpointing.

All figures show the averaged data of 5 repetitions of the benchmark runs, done at night on an idle but not dedicated cluster of different Sun workstations with SunOS 4.1.4 and a 10 Mbit/s Ethernet. Since there is no `/proc` file system, the traditional `ptrace()` interface plus the SunOS-specific extensions are used.

We used small self-made programs to measure the various aspects of interposing system calls during normal program execution, and of migration respectively checkpointing.

4.1 Costs for Interposing System Calls

There is no constant factor by which the real execution time of an application could be multiplied when running with SCIP and/or the virtual name space. The degree an application is slowed down by the interposition depends heavily on the number of system calls it does. If it does almost no system calls (e. g. only in the beginning and in the end) or if it sets itself asleep frequently, there will be almost no slow down. If it does a lot of system calls, the slow down still depends on the types of the calls.

We have measured the costs for doing system calls when the application is running in the real name space (without any system call interposition), with the modified system library ("libmod"), and with SCIP.

Calls Without Interaction with the Name Space Times for such system calls when interposed by SCIP are between 0.002 and 0.02 seconds, depending on the complexity of the interposition and the CPU speed. In all cases they are much longer than the times without interposition (factor 50...100). The corresponding times without any interposition or with interposition through the modified system library were identical and almost zero, because no interposition and no communication with the `pbeamd` is needed in these cases. However, even in the worst case the times are still short enough to be neglected, assuming a system call is a rather rare event.

One main cause for the overhead imposed by interposing via the debugging interface is the number of additional context switches between controlling and controlled process. It shows that this number is an important performance factor. So all unnecessary `ptrace()` requests should be avoided (see also below, figure 5).

Calls Requiring Interaction with the Name Space Service Now we look at some system calls that require interaction with the virtual name space administration in the `4BEAM` demon.

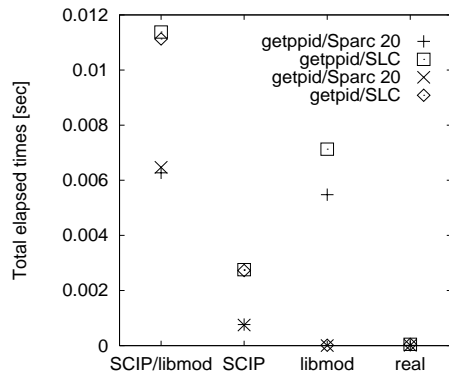


Fig. 3. Times for getting the process identifier and the parent process identifier.

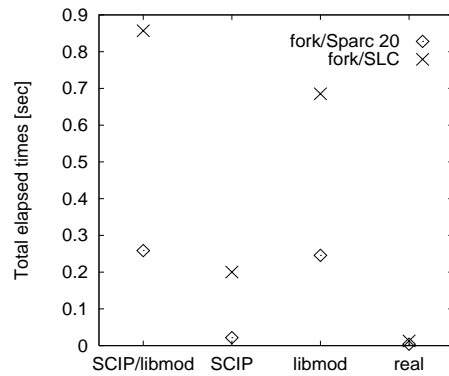


Fig. 4. Times for forking a child process.

Figure 3 shows times for querying a process' own identifier (`getpid()`), and that of its parent process (`getppid()`), measured on two machine types. The process' own virtual PID remains constant during its life time, while the parent's virtual PID may change (to 1, when the parent exits) and thus must be looked up from the name space service. The column labeled "SCIP/pbeamd" shows the times for the interposition and, if needed, lookup operation. In contrast, the column simply labeled "SCIP" shows the times for just interposing the system call *without* interaction with the pbeamd. The "libmod" and "real" curves again show the numbers for using the other respectively no interposition component.

It also shows that, because the C library functions `getpid()` and `getppid()` under SunOS are implemented using the same system call, SCIP cannot distinguish them and must always perform the expensive action of looking up the parent ID.

The most expensive system call under SCIP is creating a new process with `fork()`, since it requires forking two processes, attaching them and registering to the name space. Figure 4 shows the times needed for this, again for two different machine types and the four different kinds of interposition.

Communicating Processes To handle parallel and distributed applications, it is important to support the system calls for I/O and inter-process communication. In contrast to other system calls, whose arguments are either integers or structures of only several words, the data area to be read or written by an I/O operation is typically much larger – regarding streams, in theory, of almost arbitrary size. Nonetheless they must be considered as arguments respectively return values.

One simple approach would be to let the controlled application perform these operations, if necessary after translating any virtual to real addresses. Unfortunately, this does not work for operations that have to conform to a certain protocol in order to cope with migration and global checkpointing of distributed applications [20], i. e. especially for interprocess communication.

So another simple approach would be to let SCIP maintain these file descriptors and perform these operation on behalf of the application. This involves exchanging all the

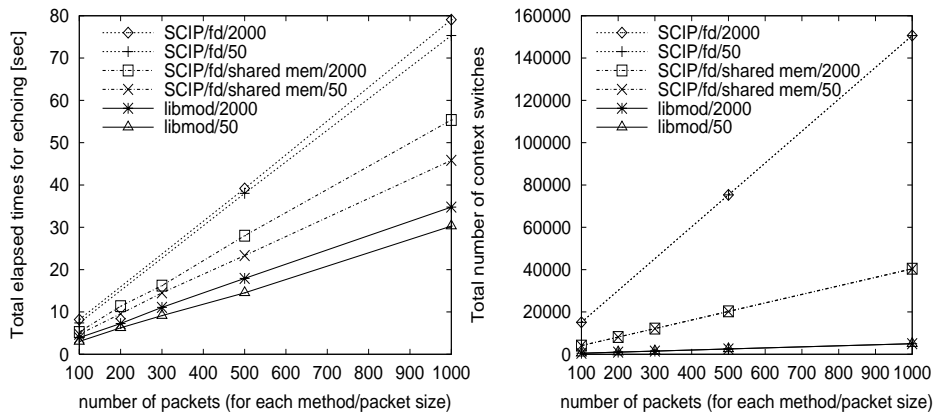


Fig. 5. Results of the pingpong benchmark: using Shared Memory saves almost half the time and most context switches.

data read or written between SCIP and the application.

Figure 5 shows times and number of context switches for a “pingpong” benchmark program. In these tests, two Sun IPCs are echoing back and forth a certain number of messages using UDP datagrams, without working on these messages. This is not a typical behavior of distributed applications, but a good example to show the cost of interposing the communication operations. The `sendto()` and `recvfrom()` operations used in the “pingpong” require to look up the real destination respectively virtual sender address in the name space service. They can send and receive messages of variable size. The figures show the benchmark results for two different packet sizes (50 and 2000 Bytes) and four different interposition methods.

Our first approach (not in the figure) was to let SCIP perform all communication on behalf of the application process. It turned out to be very expensive because the received messages must be poked word by word into pingpong’s target buffer. Now, we have implemented two optimized versions of SCIP. The first one (“SCIP/fd”) avoids copying the received message by passing the access rights to its own file descriptor to the pingpong program, that can then receive the message on its own. This involves patching and executing a short function in the application each time `recvfrom()` is called. For large packets the time savings compared to the first version are immense, but poking this function again and again onto the application’s stack is expensive, too.

Therefore, the second experimental version (“SCIP/fd/shared mem” in figure 5) puts this function into a shared memory segment, so that it can at any time be called without any overhead. This again saves most context switches and almost half the time. The measured data shows that if – and only if – such mechanisms are used where possible, the performance of SCIP can get close to that of the modified system library, that does not encounter any problems with involuntary context switches. The advantage seems worth the more difficult implementation.

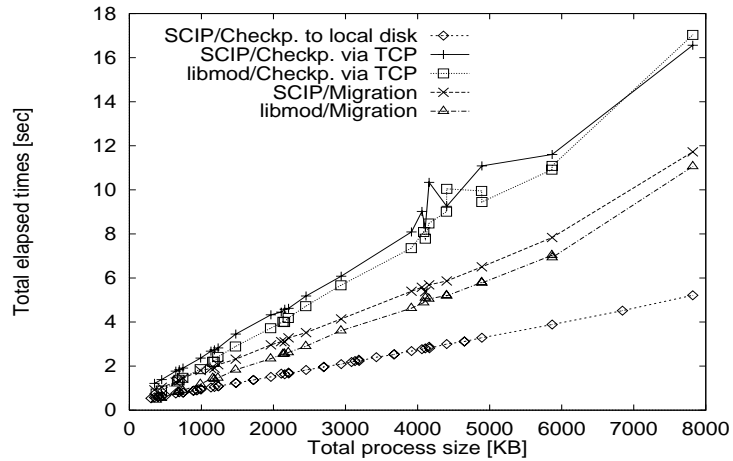


Fig. 6. Times for checkpointing/migration between two Sun IPCs.

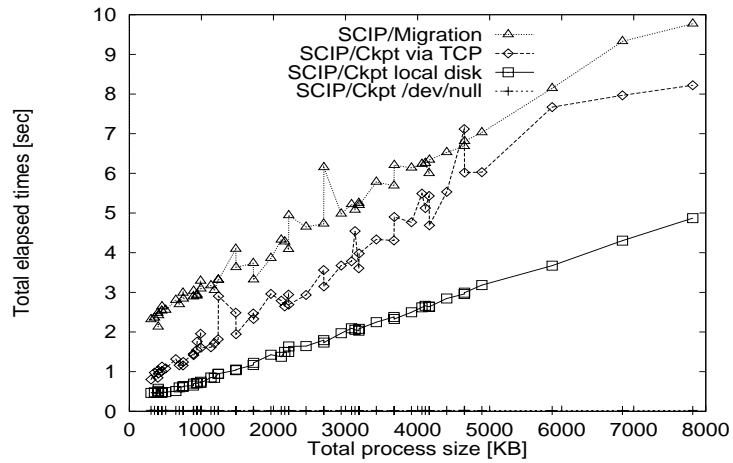


Fig. 7. Times for checkpointing/migration from a SPARCstation 10 to a SPARCstation 20.

4.2 Costs for Migration and Checkpointing

Figures 6 and 7 show the total elapsed times for checkpointing and migrating processes of varying size up to eight MB of data and stack memory. The state image is either written to a local disk file, or sent via TCP to a server process running on the target machine and written to the local disk there. This is much faster than using NFS. Migration is likewise done via TCP, as described in section 3.3.

For figure 6, both source and target machines were Sun IPC workstations. Checkpointing and migration are done using SCIP as well as the modified system library ("libmod"). As described in sections 3.3 and 3.4, the same mechanisms for reading and writing the process image are used by both methods, and the measured data confirms that both take almost the same time, with SCIP causing only a little and constant overhead. Here, the bottleneck was a slow hard disk and very low memory on the server machine.

For figure 7 the same was repeated with SCIP between an SPARCstation 10 and a SPARCstation 20. The figure shows also values for writing to the local null device for comparison. In this experiment, the bottleneck was the network bandwidth. These times, about ten seconds for eight MB process image via TCP, are only a few percent slower than with the modified system library [18], and thus SCIP still yields acceptable performance.

5 Related Work

In the last years, several process migration and checkpointing systems have been developed, which run outside the system kernel in user space only.

Condor [13, 12] provides migration and rollback transparency for sequential applications through a modified system library that forwards the system call arguments back to a shadow process on the application's home node to invoke the real service there. It does not support signals, timers or communication operations, which are needed for distributed applications.

Mandelberg and Sunderam [15] read the internal state via `ptrace()` for creating checkpoints. Information about the external state is read from kernel memory. They do not use a modified system library, thus they can handle unmodified binary programs. Applications are allowed to access files on NFS-mounted file systems only. They do not interpose system calls but collect information from kernel data structures and create links to the files by faking NFS link calls. They report that migration is slow because for one reading the address space with the "traditional" `ptrace()` facility is expensive, and because they transferred the state information via NFS. They need 3 seconds per 100 KB process image size on Sun3 machines, while SCIP achieves about 0.12...0.22 s / 100 KB.

In the ARTEMIS project [21], the operating system's shared library mechanism is exploited to link commercial applications with a modified system library. This, obviously, requires the applications to be dynamically linked, and is not applicable to statically linked programs. In Φ BEAM and SCIP we did not yet implement support for shared memory and thus not for dynamically linked applications.

Other migration systems require the applications to be based on some specific parallel programming environment like PVM [8] or MPI [16]. They provide transparency for the objects' names (e.g. task IDs) that are defined by the programming environment. They all require the applications to be linked with special versions of the parallel programming library [22, 5, 11, 25, 2]. For Φ BEAM, these programming environments are part of the application.

Trinitis is integrating an external checkpointer into CoCheck [24, 22] to do the single process checkpointing via the debugging interface. However, this is insufficient to achieve migration transparency for the message passing layer, beyond the internal state. Here, Φ BEAM/SCIP has the advantage of working at the system call level.

A very similar technique to SCIP is used by Alexandrov et al. for a different purpose [1]. Their Ufo file system extends the operating system functionality at user level with user-installable file systems. The file operation system calls are intercepted via the `/proc` file system and their arguments and return codes changed.

6 Conclusions

We conclude that the separation between the system call interposition component and the administrative components of Φ PBEAM makes it possible to choose that interposition component that is suited best for the specific task at hand. Besides the modified system library we have explored ways to transparently handle arbitrary binaries, that are already linked. The tradeoff, apart from implementation issues, is between transparency and performance. Although the approach via the debugging interface is more difficult to implement than the modified system library, we can still achieve some degree of portability with it (see also the remarks about portability of Ufo [1]).

The presented measurements show that also an acceptable performance almost similar to that of the modified system library approach can be gained. The performance is significantly better than what we initially expected from the numbers reported in the literature for previous work [15]. The main reasons for this improvement are the exploitation of the debugging interface extensions of modern Unix flavors, and even more the avoidance of as many context switches between application and interposition process as possible.

References

Many of the cited documents and some more are available online in the Internet. We have set up a World-Wide-Web page with references on <http://www.bode.informatik.tu-muenchen.de/~petri/pbeamrefs.html>.

1. A.D. Alexandrov, M. Ibel, K.E. Schauer, and C.J. Scheiman. Extending the Operating System at the User Level: the Ufo Global File System. In *USENIX Technical Conference Proceedings*, pages 77–90, Anaheim, CA, January 1997.
2. D. Andres, C. Elford, B. Fin, and L. Smith. Dynamic load balancing in PVM. Technical report, University of Illinois at Urbana-Champaign, April 1993.
3. M. Bolz. Transparent Redirection of System Calls for Unmodified Programs in Φ PBEAM. Master's thesis, Institut für Betriebssysteme und Rechnerverbund, TU Braunschweig, November 1997. (In German).
4. J. Cargille and B.P. Miller. Binary Wrapping: A Technique for Instrumenting Object Code. *ACM Sigplan Notices*, 27(6):17–18, June 1992.
5. J. Casas, D.L. Clark, R. Konuru, S.W. Otto, R.M. Prouty, and J. Walpole. MPVM: A migration transparent version of PVM. *Computing Systems*, 8(2):171–216, 1995.
6. CCS Annual Report. WWW page, Center for Computational Sciences, Oak Ridge National Laboratory, 1995. <http://www.ccs.ornl.org/AnRep95/CCS95.html>.
7. R. Faulkner and R. Gomes. The Process File System and Process Model in UNIX System V. In *USENIX Technical Conference Proceedings*, pages 243–252, Dallas, TX, January 1991.
8. Al Geist, A. Beguelin, J. Dongarra, W. Jiang, R. Manchek, and V. Sunderam. *PVM: Parallel Virtual Machine – A Users' Guide and Tutorial for Networked Parallel Computing*. The MIT Press, Cambridge, Massachusetts, 1994.
9. M.B. Jones. *Transparently Interposing User Code at the System Interface*. PhD thesis, CMU, September 1992.
10. A.H. Karp, M. Heath, and Al Geist. 1995 Gordon Bell Prize Winners. *IEEE Computer*, 29(1):79–85, January 1996.

11. J. León, A.L. Fisher, and P. Steenkiste. Fail-save PVM: A portable package for distributed programming with Transparent Recovery. Report CMU-CS-93-124, Carnegie Mellon University, February 1993.
12. M. Litzkow, T. Tannenbaum, J. Basney, and M. Livny. Checkpointing and Migration of UNIX Processes in the Condor Distributed Processing System. Report 1346, University of Wisconsin-Madison Computer Sciences, April 1997.
13. M.J. Litzkow and M. Solomon. Supporting Checkpointing and Process Migration Outside the UNIX Kernel. In *USENIX Technical Conference Proceedings*, pages 283–290, San Francisco, CA, January 1992.
14. D. Long, J. Carroll, and C. Park. A Study of the Reliability of Internet Sites. In *Proceedings of the 10th Symposium on Reliable Distributed Systems*, pages 177–186, 1991.
15. K.I. Mandelberg and V.S. Sunderam. Process Migration in UNIX Networks. In *USENIX Technical Conference Proceedings*, pages 357–363, Dallas, TX, February 1988.
16. Message Passing Interface Forum MPIF. MPI-2: Extensions to the Message-Passing Interface. Technical report, University of Tennessee, Knoxville, July 1997. <http://www.mpi-forum.org>.
17. S. Petri, M. Bolz, and H. Langendörfer. Transparent Migration and Rollback for Unmodified Applications in Workstation Clusters. Informatik-Bericht 98-02, TU Braunschweig, April 1998. To appear.
18. S. Petri and H. Langendörfer. Load Balancing and Fault Tolerance in Workstation Clusters – Migrating Groups of Communicating Processes. *Operating Systems Review*, 29(4):25–36, October 1995.
19. S. Petri, B. Schnor, M. Becker, B. Hinrichs, T. Tschardtke, and H. Langendörfer. Evaluation of Multicast Methods to Maintain a Global Name Space for Transparent Process Migration in Workstation Clusters. In *Kommunikation in Verteilten Systemen*, pages 224–234. GI/ITG Fachtagung KIVS'97, Springer, February 1997.
20. S. Petri, B. Schnor, H. Langendörfer, and J. Steinborn. Consistent Global Checkpoints for Distributed Applications on Clusters of Unix Workstations. In *Paralleles und Verteiltes Rechnen – Beiträge zum 4. Workshop über Wissenschaftliches Rechnen*, pages 77–86, Aachen, October 1996. TU Braunschweig, Shaker.
21. T. Shirakihara, H. Hirayama, K. Sato, and T. Kanai. ARTEMIS: Advanced Reliable distributed Environment Middleware System. In *Proceedings of the International Conference on Parallel and Distributed Processing Techniques and Applications, PDPTA'97*, pages 97–106, Las Vegas, NV, July 1997.
22. G. Stellner. CoCheck: Checkpointing and Process Migration for MPI. In *Proceedings of the 10th International Parallel Processing Symposium (IPPS '96)*, Honolulu, Hawaii, April 1996.
23. Sun Microsystems. *SunOS Reference Manual*, 1990. Revision A.
24. J. Trinitis. An External Checkpointing Technique for Integration into a Parallel Tool Environment. In preparation. trinitis@informatik.tu-muenchen.de, 1998.
25. J.J.J. Vesseur, R.N. Heederik, B.J. Overeinder, and P.M.A. Sloot. Experiments in Dynamic Load Balancing for Parallel Cluster Computing. In *Proceedings of the Workshop on Parallel Programming and Computation (ZEUS'95) and the 4th Nordic Transputer Conference (NTUG'95)*, pages 189–194, Amsterdam, June 1995. IOS Press.