

VU Research Portal

Amoeba on a Multiprocessor

van Moergestel, L.J.; Bal, H.E.; Kaashoek, M.F.; van Renesse, R.; Sharp, G.J.

1989

document version

Publisher's PDF, also known as Version of record

[Link to publication in VU Research Portal](#)

citation for published version (APA)

van Moergestel, L. J., Bal, H. E., Kaashoek, M. F., van Renesse, R., & Sharp, G. J. (1989). *Amoeba on a Multiprocessor*. (Report, Dept. of Mathematics and Computer Science; No. IR-206).

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal ?

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

E-mail address:

vuresearchportal.ub@vu.nl

Amoeba on a Multiprocessor

Leo J. M. van Moergestel

Henri E. Bal¹

Frans Kaashoek

Robbert van Renesse¹

Gregory J. Sharp

Hans van Staveren

Andrew S. Tanenbaum

Dept. of Mathematics and Computer Science

Vrije Universiteit

De Boelelaan 1081

1081 HV Amsterdam, The Netherlands

Internet: leovm@cs.vu.nl or amoeba@cs.vu.nl

ABSTRACT

The Amoeba distributed operating system has been in development and use for over eight years now. In this paper we describe the hardware of the processor pool on which Amoeba runs.

Computing Reviews categories: C.2.4, D.4

Keywords: Multiprocessors, VMEbus, Operating systems, Distributed systems, Distributed operating systems

1. Introduction

The Amoeba project is a research effort aimed at understanding how to connect multiple computers together in a seamless way (Mullender and Tanenbaum, 1986; Tanenbaum et al., 1986; Tanenbaum and van Renesse, 1985). The basic idea is to provide the users with the illusion of a single powerful timesharing system, when, in fact, the system is implemented on a collection of machines, potentially close together, in a single rack in the machine room, or far apart, in different countries. This research has led to the design and implementation of the Amoeba distributed operating system, which is being used as a prototype and vehicle for further research. In this paper we will describe some hardware aspects as well as an application using a language, named *Orca*, for parallel applications running on a multiprocessor.

Amoeba was originally designed and implemented at the Vrije Universiteit in

1. This research was supported in part by the Netherlands Organization for Scientific Research (N.W.O.) under grant 125-30-10.

Amsterdam, and is now being jointly developed there and at the Centre for Mathematics and Computer Science, also in Amsterdam. The chief goal of this work is to build a distributed system that is *transparent* to the users. This concept can best be illustrated by contrasting it with a network operating system, in which each machine retains its own identity. With a network operating system, each user logs into one specific machine, his home machine. When a program is started, it executes on the home machine, unless the user gives an explicit command to run it elsewhere. Similarly, files are local unless a remote file system is explicitly mounted or files are explicitly copied. In short, the user is clearly aware that multiple independent computers exist, and must deal with them explicitly.

In a transparent distributed system, in contrast, users effectively log into the system as a whole, and not to any specific machine. When a program is run, the system, not the user, decides the best place to run it. The user is not even aware of this choice. Finally, there is a single, system-wide file system. The files in a single directory may be located on different machines possibly in different countries. There is no concept of file transfer, uploading or downloading from servers, or mounting remote file systems. The fact that a file is remote is not visible to the user at all.

The remainder of this paper will give a short overview of Amoeba, a description of the hardware and some speed measurements of an application running on a multiprocessor.

2. Amoeba System Architecture

The Amoeba architecture consists of four principal components, as shown in Fig. 1. First are the workstations, one per user, on which users can carry out editing and other tasks that require fast interactive response. The workstations are all diskless, and are primarily used as intelligent terminals that do window management, rather than as computers for running complex user programs. We are currently using SUN-3s and VAXstations as workstations. In the next generation of hardware we may use X-terminals.

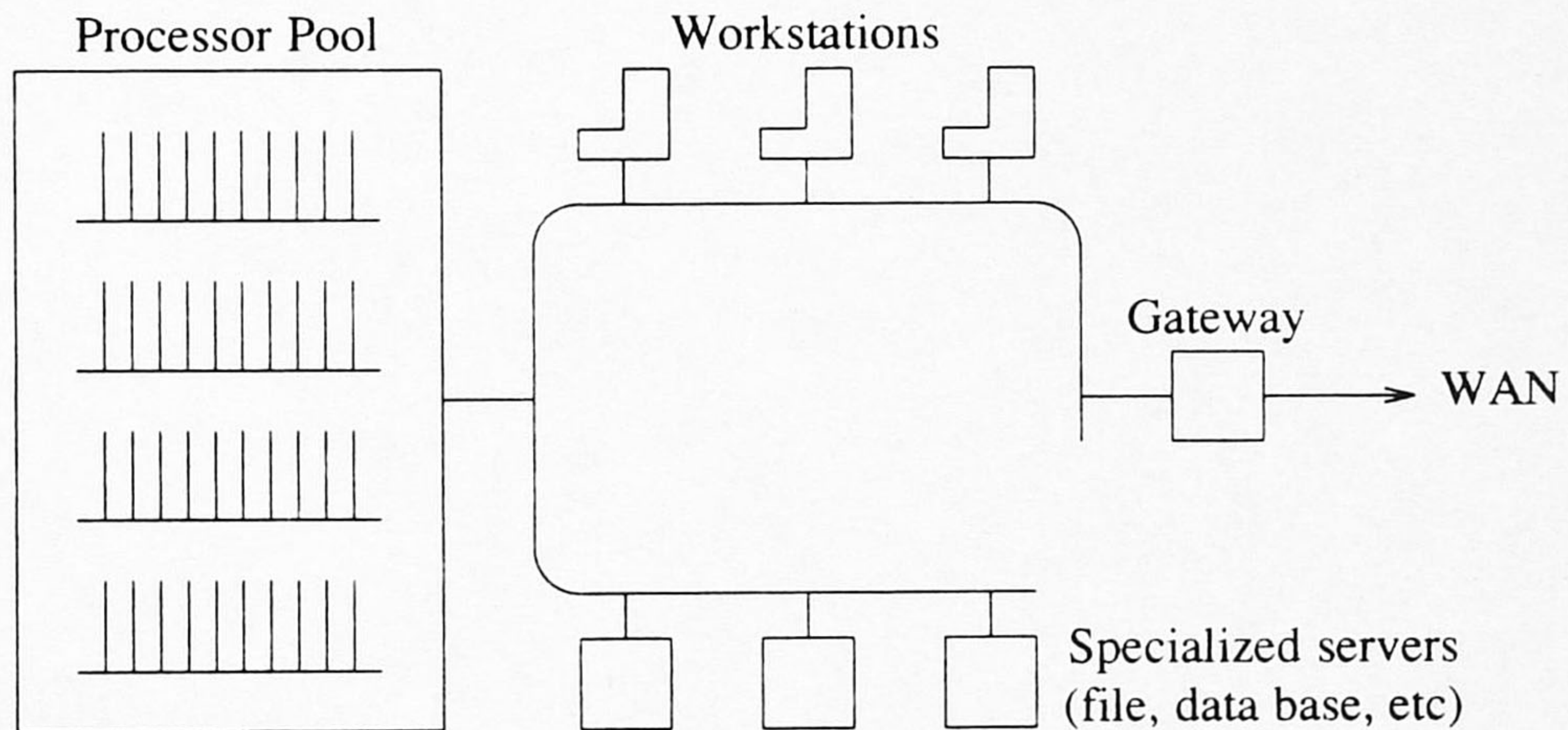


Fig. 1. The Amoeba architecture.

Second are the pool processors, a group of CPUs that can be dynamically allocated as

needed, used, and then returned to the pool. For example, the *make* command might need to do six compilations, so six processors could be taken out of the pool for the time necessary to do the compilation and then returned. Alternatively, with a five-pass compiler, 5 times 6 = 30 processors could be allocated for the six compilations, gaining even more speedup. Many applications, such as heuristic search in AI applications (e.g., playing chess), use large numbers of pool processors to do their computing. At the Vrije Universiteit we currently have 42 single board VME-based computers using the MC68020 and MC68030 CPUs.

Third are the specialized servers, such as directory servers, file servers, data base servers, boot servers, and various other servers with specialized functions. Each server is dedicated to performing a specific function. In some cases, there are multiple servers that provide the same function, for example, as part of the replicated file system.

Fourth are the gateways, which are used to link Amoeba systems at different sites and different countries into a single, uniform system. The gateways isolate Amoeba from the peculiarities of the protocols that must be used over the wide-area networks.

All the Amoeba machines run the same kernel, which primarily provides multithreaded processes, communication services, and little else. The basic idea behind the kernel was to keep it small, to enhance its reliability, and to allow as much as possible of the operating system to run as user processes, providing for flexibility and experimentation.

3. A Closer Look at the Processor Pool

Communication between CPUs is the basic hardware requirement for the Amoeba processor pool.

3.1. LAN coupled CPU boards

The simplest setup from a hardware point of view is using stand alone processors connected by a Local Area Network (LAN). We have an experimental setup with 16 MC68020 boards with 2Mbyte RAM each connected to a 10Mbit/sec Ethernet. We call this a *loosely coupled* system, because the processors only communicate over Ethernet; there is no shared memory. The advantage is the great flexibility: you just tie more CPUs to the LAN to expand your computing capacity. You can do this while Amoeba is alive and running on the other systems. The drawback is that the bandwidth of the LAN poses a severe limit on your communication speed. When the LAN is being used for other purposes, even less bandwidth is available.

3.2. Multiprocessor

For applications requiring a lot of data communication, one can use a higher bandwidth medium. This results in a *tightly coupled* system in which the processors share the same back-plane and bus.

We chose the VME bus as a universal board level bus for communication. Currently we are using two multiprocessor systems based on the VME bus. One is built with 10 MC68020 CPU boards. Another one contains 17 MC68030 boards. We will give a more detailed description of the latter. It consists of 16 CPU30ZA boards and a CPU30ZAE with Ethernet interface to talk to the outside world. In this paper we will call the board with Ethernet interface the *master*. The other boards are the *slaves*. This has nothing to do with the VME concept

of master and slave. An ASCU-2 bus controller and a VME RAM board are also part of the system. The CPU boards are manufactured by Force and contain a MC68030 CPU with 4Mbyte dual ported RAM. Furthermore the on-board FGA002 programmable gate array gives a good support for multiprocessing.

For the system to function we had to solve the following problems.

- 1) All boards should be able to access the bus even when it is heavily loaded.
- 2) Each board should be able to communicate with all the others.

3.2.1. Arbitration.

The VME bus has four levels for bus arbitration. The prioritised arbitration scheme would be a bad choice, because the higher-level boards could take all the bus bandwidth leaving nothing over for the lower level boards. Moreover, our multiprocessor concept is based on the idea that all processors have equal access to the bus as all CPUs have equal access to the LAN in the loosely coupled system. The VME specification also offers a round robin arbitration, which is better, but within a level we keep a prioritised daisy chain for possible bus masters with decreasing priority along the chain. The FGA002 solves this problem by introducing a fair arbitration scheme within a certain arbitration level. This means when a board gets the bus due to the round robin arbiter giving the level for arbitration it will not take the bus again if there is a lower level board in the daisy chain asking for the bus. This results in fair bus access for the lower boards in the daisy chain even in a heavy loaded system.

To demonstrate the effect of fair arbitration we wrote a test program. The master downloads the slaves and generates a broadcast message (Force Message Broadcast, FMB) to start the test. During the test the slaves try to increment a counter in global (VME) memory. A reset will stop the process and the counter values can be read.

We ran this test with 16 boards (16 slaves) Each bus request level had 4 boards in a daisy chain. Processor 1, 2, 3 and 4 were on BREQ 3; processors 5, 6, 7, and 8 were on BREQ 2 and so on. A FMB interrupt set a flag to TRUE to start the counting. Interrupts were off during counting. Fair arbitration was on in one set of tests and off in another. Between two increment instructions the CPUs entered a small delay-loop.

MASTER program
initialise counters to zero.
generate FMB interrupt.

SLAVE program
initialise FMB interrupt.
get counter address by using own VME slot number as an index in the array of counters.
wait for flag to become TRUE.
disable interrupts.
start incrementing counter.

Interrupt handler in SLAVE
FMB interrupt: make startflag TRUE.

Figure 2 and 3 show the result for a delay of 10 arbitrary units. In these figures we see that in the normal arbitration scheme, number four in the daisy chain is completely unable to access the bus. In fact this CPU encountered a bus timeout error so it gave simply up. When we just skip the delay between the increments the results are even more dramatic as figure 4 and 5 show.

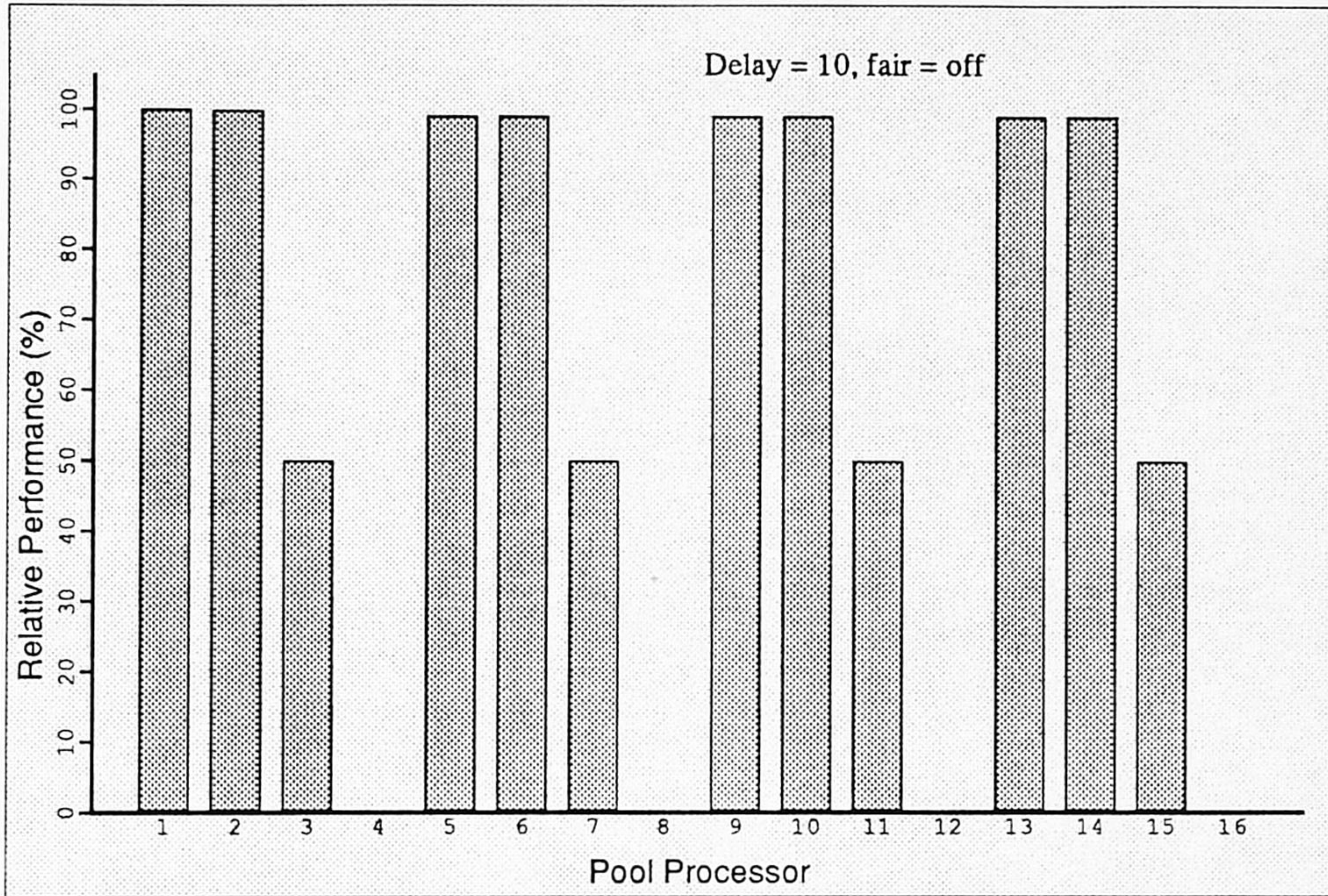


Fig. 2. Fair arbitration off

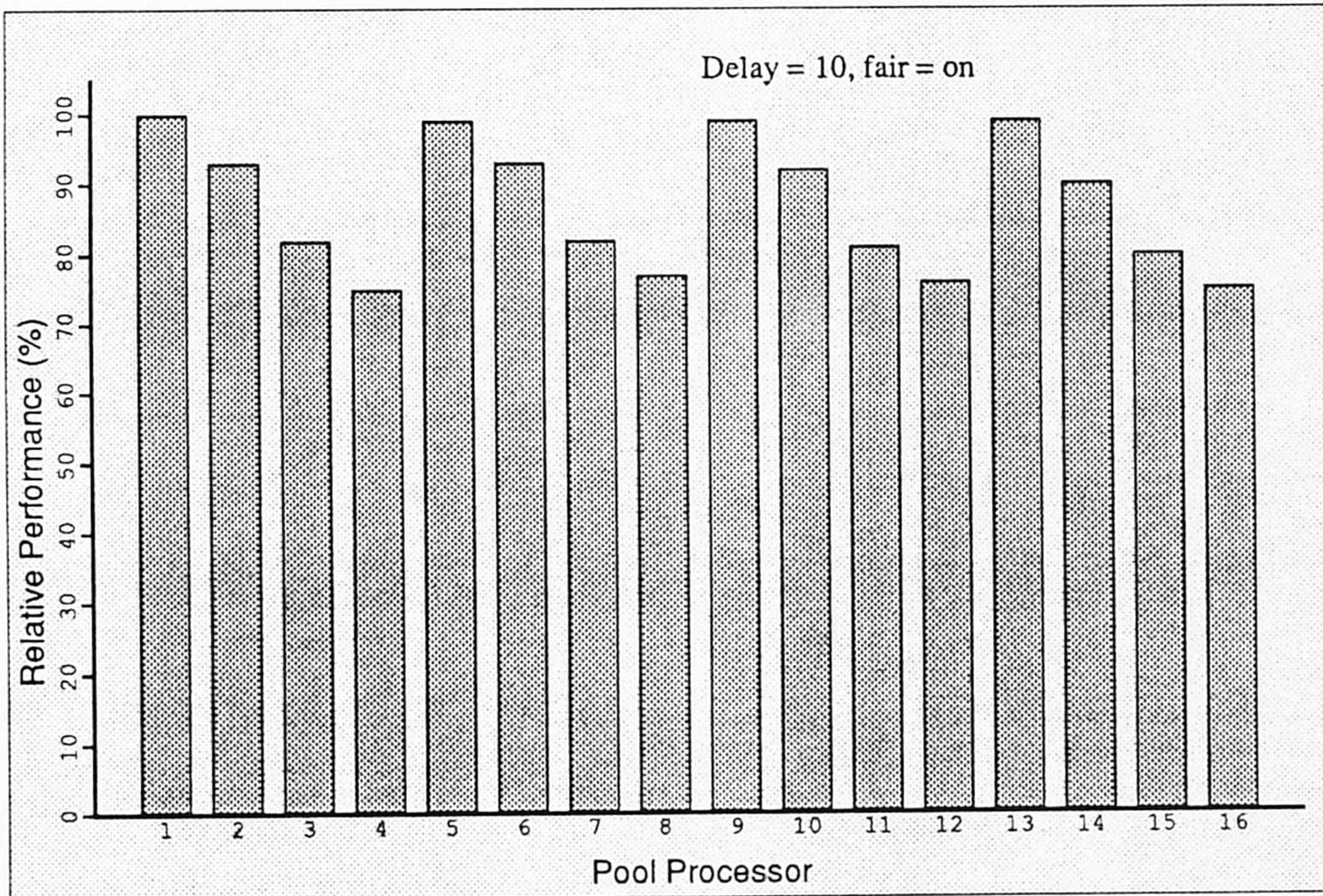


Fig. 3. Fair arbitration on

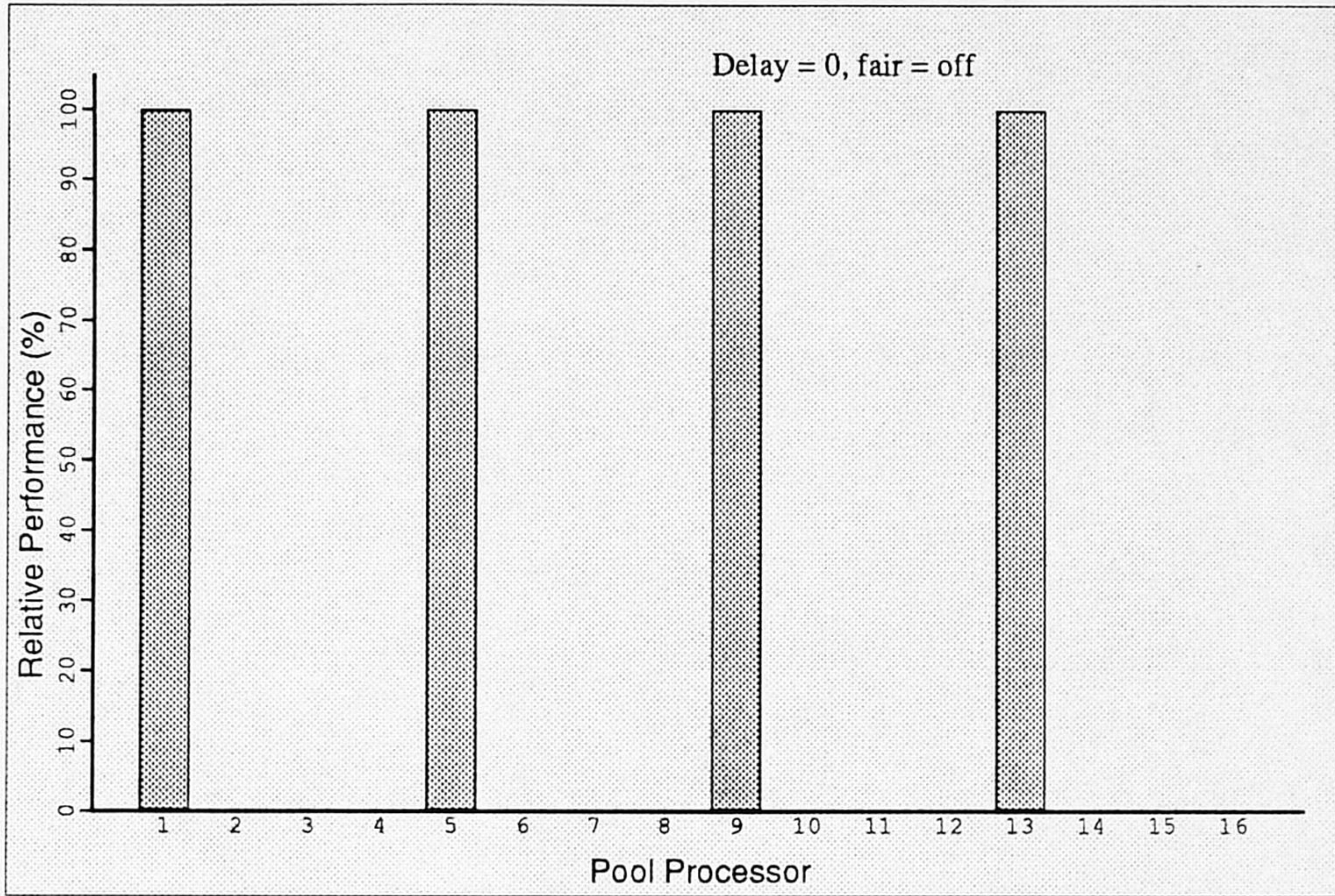


Fig. 4. Fair arbitration off

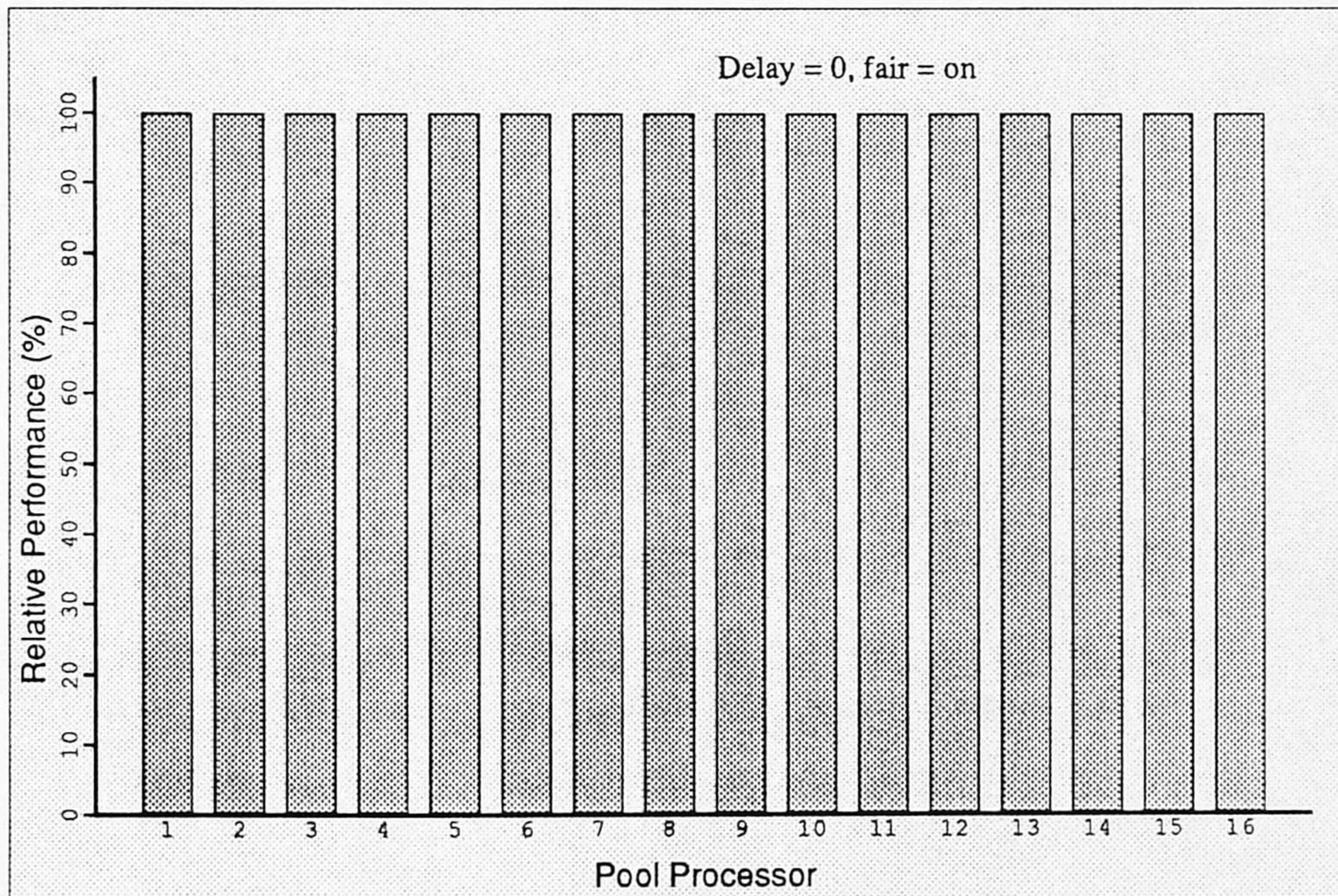


Fig. 5. Fair arbitration on

3.2.2. Interboard Communication Using Location Monitors.

The way the CPUs communicate is by writing to a certain buffer in the destination processor's dual ported RAM and then generating a mailbox interrupt. Each processor 'owns' a buffer in the dual ported RAM of the other processor boards. The simplest way to implement a location monitor or mailbox is by generating an interrupt when there is an access to a certain address. This implementation has a drawback. After generating the interrupt, the interrupted CPU has to handle it, which takes a certain amount of time. If another board also decides to generate a mailbox interrupt simultaneously, one of the interrupts will be lost, depending on the software implementation of the interrupt handling. Of course one can solve this problem, called a *race condition*, in software, but it would be nice if the hardware helped.

The implementation of the location monitors on the CPU30 is more like a semaphore. To generate an interrupt, the interrupting CPU reads a byte from a location monitor address. If the read returns a zero the mailbox was free and an interrupt is generated on the CPU accessed. Until the interrupted CPU does a write to the same mailbox, all other boards that try to use this mailbox will get a non-zero result on a read, meaning that the mailbox is not empty and no interrupt is being generated. They thus 'know' they did not get through and have to try again or use another mailbox. The interrupted CPU will eventually free the mailbox by writing to it (in the interrupt service routine).

3.2.3. Dual Ported RAM

The on-board RAM is accessible by both the CPU and the VME bus. This means that all CPUs are able to access each other's on board RAM. With the help of the FGA002 it's possible to protect some parts of this RAM so one board with corrupted code cannot destroy the programs residing on the other boards. Only the communication buffers are accessible.

This brings us to another topic. What if one board crashes? The master board discovers that it cannot communicate anymore with the crashed board. Again the gate array offers a solution. By writing to a certain register (accessible from the VME bus) an on board RESET is generated. So the failing board can be rebooted without disturbing the whole system.

4. Low Level Boot Software.

After initialising the gate array. A small program in ROM is started on the slave boards. This program waits for a mailbox interrupt. When it occurs, a certain RAM location is checked for a magic number and a jump to that location is executed.

The master board containing the Ethernet interface runs a small Amoeba kernel in ROM so it can communicate with the outside world. The master is downloaded with the Amoeba kernel for the slaves as well with its own runtime system. It examines the VME bus and copies the slave kernel to all the slaves it discovers. The slaves are then triggered by a mailbox interrupt and start executing their code.

5. Applications

Amoeba has been used to program a variety of applications. To mention a few: UNIX†

† UNIX is a Registered Trademark of AT&T Bell Laboratories.

emulation, parallel *make*, traveling salesman, and alpha-beta search.

5.1. Parallel Applications and Orca

Although Amoeba was originally conceived as a system for *distributed* computing, the existence of the multiprocessor with 16 MC68030 CPUs close together has made it quite suitable for *parallel* computing as well. That is, we have become much more interested in using the multiprocessor to achieve large speedups on a single problem. To program these parallel applications, we have designed and implemented a language called Orca (Bal and Tanenbaum, 1988).

Orca is based on the concept of globally shared objects. Programmers can define operations on shared objects, and the compiler and run time system take care of all the details of making sure they are carried out correctly. This scheme gives the programmer the ability to atomically read and write shared objects that are physically distributed among a collection of machines without having to deal with any of the complexity of the physical distribution. All the details of the physical distribution are completely hidden from the programmer. Initial results indicate that almost linear speedup can be achieved on some problems.

5.2. Example of Problem Solving Using Orca

In this section we will describe traveling salesman problem. In the traveling salesman problem, the computer is given a starting city and a list of other cities to be visited. The idea is to find the shortest path that visits each city exactly once, and then returns to the starting place. Using Amoeba we have programmed this application in parallel by having one pool processor act as coordinator, and the rest as slaves.

Suppose, for example, that the starting place is London, and the cities to be visited include New York, Sydney, Nairobi, and Tokyo. The coordinator might tell the first slave to investigate all paths starting with London-New York, the second slave to investigate all paths starting with London-Sydney, the third slave to investigate all paths starting with London-Nairobi, and so on. All of these searches go on in parallel. When a slave is finished, it reports back to the coordinator and gets a new assignment.

The algorithm can be applied recursively. For example, the first slave could allocate a processor to investigate paths starting with London-New York-Sydney, another processor to investigate London-New York-Nairobi, and so forth. At some point, of course, a cutoff is needed at which a slave actually does the calculation itself and does not try to farm it out to other processors.

The performance of the algorithm can be greatly improved by keeping track of the best total path found so far. A good initial path can be found by using the "closest city next" heuristic. Whenever a slave is started up, it is given the length of the best total path so far. If it ever finds itself working on a partial path that is longer than the best-known total path, it immediately stops what it is doing, reports back failure, and asks for more work.

Figure 6 shows the results of a timing measurement of the Traveling salesman problem as a function of the total number of CPUs used to perform the calculation.

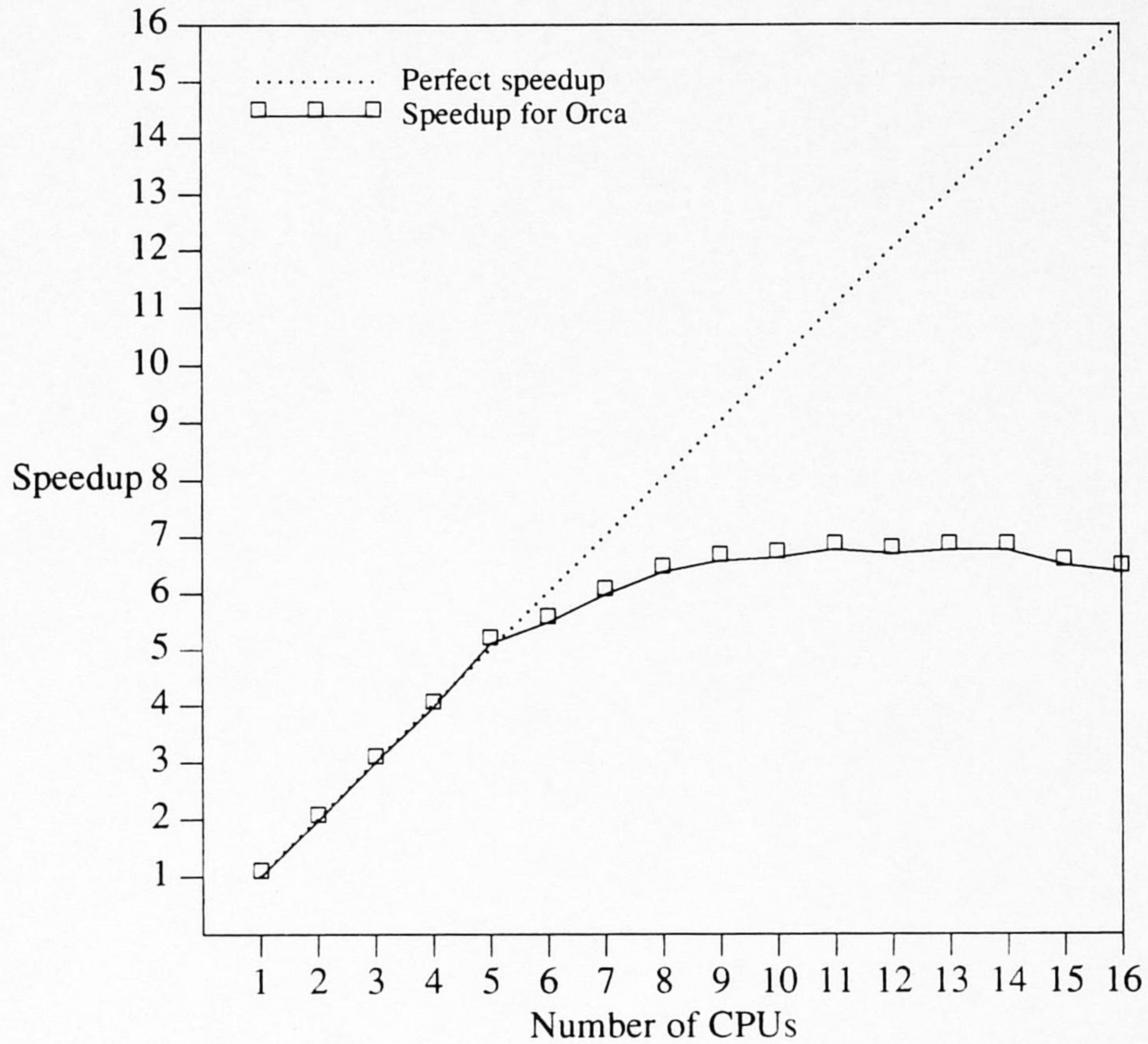


Fig. 6. Traveling salesman

This picture shows that we have a problem when too many CPUs try to read the 'best total path found so far' in global VME RAM. This problem can be solved by keeping a copy of this value in local RAM and perform an update whenever a better result is found.

REFERENCES

- Bal, H.E., Renesse, R. van, and Tanenbaum, A.S. Implementing Distributed Algorithms using Remote Procedure Call, *Proc. National Computer Conference*, AFIPS, 1987. pp. 499-505.
- Bal, H.E., and Tanenbaum, A.S. Distributed Programming with Shared Data, *IEEE Conf. on Computer Languages*, IEEE, 1988, pp. 82-91.
- Cheriton, D.R. The V Distributed System, *Commun. ACM* 31, (March 1988), pp. 314-333.
- Marsland, T.A., and Campbell, M. Parallel Search of Strongly Ordered Game Trees, *Computing Surveys* 14, (Dec. 1982) pp. 533-551.

- Mullender, S.J., and Tanenbaum, A.S. The Design of a Capability-Based Distributed Operating System, *Computer Journal* 29, (Aug. 1986), pp. 289-299.
- Mullender, S.J., and Tanenbaum, A.S. A Distributed File Service Based on Optimistic Concurrency Control, *Proc. Tenth Symp. Oper. Syst. Prin.*, 1985, pp. 51-62.
- Renesse, R. van, Tanenbaum, A.S., and Wilschut, A The Design of a High-Performance File Server *Proc. Ninth Int'l Conf. on Distr. Comp. Systems*, IEEE, 1989a, pp. 22-27.
- Renesse, R. van, Staveren, H. van, and Tanenbaum, A.S. Performance of the Amoeba Distributed Operating System, *Software—Practice and Experience* 19, (March 1989b) pp. 223-234.
- Renesse, R. van, Staveren, H. van, and Tanenbaum, A.S. Performance of the World's Fastest Distributed Operating System, *Operating Systems Review* 22, (Oct. 1988), pp. 25-34.
- Renesse, R. van, and Tanenbaum, A.S. Voting with Ghosts, *Proc. Eighth Int'l Conf. on Distr. Comp. Systems*, IEEE, 1988, pp. 456-461.
- Renesse, R. van, Tanenbaum, A.S., Staveren, H. van, and Hall, J. Connecting RPC-Based Distributed Systems Using Wide-Area Networks, *Proc. Seventh Int'l Conf. on Distr. Comp. Systems*, IEEE, 1987, pp. 28-34.
- Tanenbaum, A.S. A UNIX Clone with Source Code for Operating Systems Courses, *Operating Syst. Rev.* 21, (Jan. 1987), pp. 20-29.
- Tanenbaum, A.S., Mullender, S.J., and Renesse, R., van Using Sparse Capabilities in a Distributed Operating System *Proc. Sixth International Conf. on Distr. Computer Systems*, IEEE, 1986.
- Tanenbaum, A.S., and Renesse, R. van A Critique of the Remote Procedure Call Paradigm *Proc. Euteco '88* 1988, pp. 775-783.
- Tanenbaum, A.S., and Renesse, R. van Distributed Operating Systems, *Computing Surveys* 17, (Dec. 1985) pp. 419-470.