

Evaluation of Cluster Interconnects for a Distributed Shared Memory *

Sumit Roy and Vipin Chaudhary
Parallel and Distributed Computing Laboratory
Wayne State University
Detroit, MI 48202
(sroy | vchaud)@ece.eng.wayne.edu

Abstract

Clusters of Symmetrical Multiprocessors (SMPs) have recently become popular as low-cost, high-performance computing solutions. The type of interconnection hardware used in these clusters can become a deciding factor in their overall performance. This paper evaluates the performance of three different communication systems, 10 Mbps Ethernet, 100 Mbps FastEthernet and 155 Mbps ATM, using a multithreaded Distributed Shared Memory system, Strings.

The raw performance of each network is first measured using netperf. Ten different applications are then used for performance evaluation, including programs from the SPLASH-2 benchmarks, a medical computing application, and some computational kernels. It is found that half of the programs tested are not significantly affected by changes in the bandwidth. Though the ATM network provides the highest overall bandwidth, the remaining applications show that the increase in latency compared to FastEthernet prevents any performance improvement. On the other hand, applications that require only moderately high bandwidths perform substantially better with FastEthernet.

1 Introduction

Though current microprocessors are getting faster at a very rapid rate, there are still some very large and complex problems that can only be solved by using multiple cooperating processors. These problems include the so-called *Grand Challenge Problems*, such as Fuel combustion, Ocean modeling, Image understanding, and Rational drug design. Recently many vendors of traditional workstations have adopted a design strategy wherein multiple state-of-the-art microprocessors are used to build high performance shared-memory parallel workstations. These symmetrical multiprocessors (SMPs) are then connected through high speed networks or switches to form scal-

able computing clusters. This important class of machines includes the SGI Power Challenge Array using a HIPPI interconnect, the IBM RS/6000 SP2 series with multiple PowerPC based nodes, the Convex Exemplar series of machines, the DEC AdvantageCluster 5000 with the GIGAswitch, and the Cray/SGI Origin 2000 series with the scalable Craylink network.

Using multiple nodes on such SMP clusters requires the programmer to either write explicit message passing programs, using libraries like MPI or PVM; or to rewrite the code using a new language with parallel constructs eg. HPPF or Fortran 90. Message passing programs are cumbersome to write and may have to be tuned for each individual architecture to get the best possible performance. Parallel languages only work well with code that has regular data access patterns. In both cases the programmer has to be intimately familiar with the application program as well as the target architecture. The shared memory model on the other hand, is easier to program since the programmer does not have to worry about the data layout and does not have to explicitly send data from one process to another. Hence, an alternate approach to using these computing clusters is to provide an illusion of logically shared memory over physically distributed memory, known as a Distributed Shared Memory (DSM) or Shared Virtual Memory (SVM). Recent research projects with DSMs have shown good performance, for example TreadMarks [1], Quarks [2], CVM [3], CASHMERE-2L [4], and *Strings* [5].

This paper evaluates the performance of a multithreaded DSM, *Strings* using three different networking technologies, 10 Mbps Ethernet, 100 Mbps FastEthernet and 155 Mbps OC-3 ATM. The performance is evaluated using programs from the SPLASH-2 [6] benchmark suite, an application from medical computing, and some computational kernels. It is found that 50 % of the programs tested do not require a very high bandwidth, and show only moderate performance differences when using the various networks. Though some of the remaining applications are latency tolerant, the performance improvement with ATM is not very high as compared to FastEthernet. Thus the

*This research was supported in part by NSF grants MIP-9309489, EIA-9729828, US Army Contract DAEA 32-93D004 and Ford Motor Company grants 96-136R and 96-628R

added cost of this network technology does not translate into corresponding gains in program execution time. On the other hand, some applications do not require very high bandwidth and perform substantially better with FastEthernet.

The rest of the paper is organized as follows: Section 2 provides background on the distributed shared memory system used for the evaluation. Section 3 describes the experimental environment. Section 4 provides details about the programs used for testing. Section 5 presents and analyzes the results from the experiments. Section 6 summarizes this paper and suggests directions for future work.

2 Strings Distributed Shared Memory System

Strings [5] is a fully multithreaded DSM. Its distinguishing feature is that it incorporates Posix1.c threads multiplexed on kernel light-weight processes for better performance. The kernel can schedule independent threads across multiple processors on SMP nodes, using these lightweight processes. Thus, *Strings* is designed to exploit data parallelism at the application level and task parallelism at the DSM system level. Since DSMs share data at the relative large granularity of a page, it can happen that two processes try to write to different parts of a page at the same time. This phenomenon is known as false sharing and can lead to excessive ping-ponging of the page if only a single writer is allowed. *Strings* solves this problem by allowing multiple writers, and implements data coherency using a relaxed consistency model [7].

Early generation DSM systems used interrupt driven I/O to obtain pages, locks etc. from remote nodes. This can cause considerable disruption at the receiving node, and previous research tried to overcome this by aggregating messages, reducing communication by combining synchronization with data, and other such techniques [8]. *Strings* avoids this overhead by using a dedicated communication thread, which monitors the network port. Incoming message queues are maintained for each active thread at a node, and message arrival is announced using condition variables. This prevents wasting CPU cycles with busy waits. A reliable messaging system is implemented on top of UDP.

Table 1 shows the packet sizes that occur for common events associated with the *Strings* runtime. The page size for the evaluation environment was 8192 bytes. The system maintains a distributed locking scheme, and the current list of requests is sent to the next acquirer of a lock. Hence the lock grant message can contain a variable number of bytes depending on the queue length. The actual average packet size for an application depends on the relative proportion of paging, locking and barrier events.

Event	Request (bytes)	Reply (bytes)
Get Page	40	8226
Lock Acquire	36	28
Lock Grant	≥ 36	28
Barrier Arrival	32	28
Barrier Release	32	28

Table 1: Typical Message Sizes in *Strings*

3 Experimental Environment

All the experiments were carried out on a cluster of four SUN UltraEnterprise Servers. One machine is a six processor UltraEnterprise 4000 with 1.5 Gbyte memory. The master process for *Strings* was always run on this machine. The other machines are four processor UltraEnterprise 3000s, with 0.5 Gbyte memory each. All machines use 250 MHz UltraSparcII processors, with 4 Mb external cache. Three different networks are used to interconnect the machines in the cluster: 10 Mbps Ethernet with a generic hub, 100 Mbps FastEthernet with a BayStack FastEthernet Hub, and 155 Mbps OC-3 ATM with a Fore-RunnerLE 155 ATM switch,

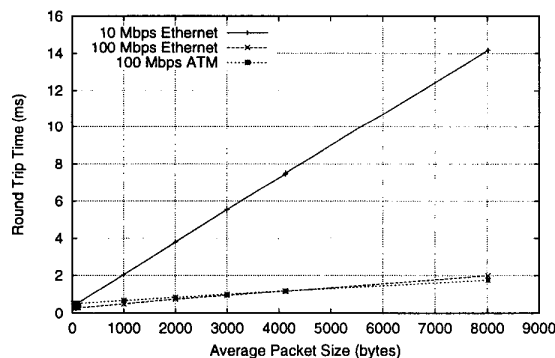


Figure 1: Roundtrip times for UDP packets

The raw roundtrip latencies and streaming bandwidth that can be obtained for UDP packets are shown in Figures 1 and 2. The results were obtained using two nodes on an unloaded network, using the *netperf-2.1* evaluation program [9]. These values place an upper bound on the performance that each network can provide, since they do not include the effects of network collisions and congestion at the receiving ports. From Figure 1, the latencies for the FastEthernet and ATM network indicate a crossover point at around 4000 bytes. Figure 2 shows that both 10 Mbps Ethernet and FastEthernet can utilize the full bandwidth for packets larger than 512 bytes. The ATM network has

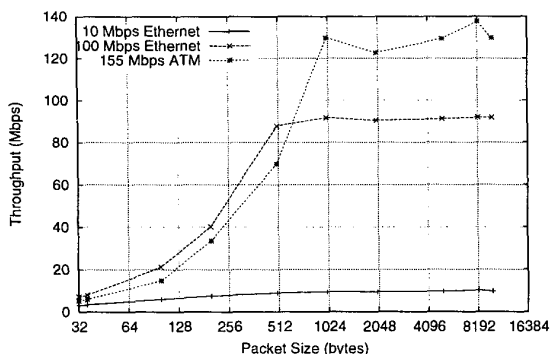


Figure 2: Throughput for UDP packets

a high startup latency, since the bandwidth utilized is less than for FastEthernet for packet sizes less than 512 bytes. The ATM network uses the full bandwidth for packets of approximately 1024 bytes.

Previous work on DSMs has presented performance on ATM [10, 1], 100 Mbps Switched FastEthernet [11], and the IBM SP-2 switch [12, 13].

4 Test Programs

The suite of test programs used to do the evaluation consists of some programs from the SPLASH-2 benchmark suite [6], a complete application program for deblurring images obtained from Magnetic Resonance Imaging (MRI), matrix multiplication (MATMULT), as well as a kernel for doing Successive Over Relaxation program (SOR) on a rectangular grid of numbers.

4.1 SPLASH-2 Benchmark Programs

The SPLASH-2 Benchmark programs have been written for evaluating the performance of shared address-space multiprocessors and include application kernels as well as full fledged code. These programs can be adapted for use on a DSM, since the shared memory initialization, synchronization primitives, and task creation are encapsulated within PARMACS-ANL macros [14]. Hence they are often used as benchmarks when presenting results on DSM systems.

The execution model for SPLASH-2 programs follows the Single Program Multiple Data (SPMD) type, and has three phases. In the first phase the main task reads in command-line parameters and the initial data. It sets up a shared memory region and allocates globally shared data from this region. Synchronization variables are also initialized. After this initialization phase, tasks are created to execute the actual *slave()* routine in the computation phase. The master also runs a copy of this routine. In the termination phase the master thread collects the results and gener-

ates timing statistics.

The data access patterns of the programs in the SPLASH-2 suite have been characterized in earlier research [15, 16]. FFT performs a one-dimensional Fast Fourier Transform of n complex data points. Three all-to-all interprocessor communication phases are required for a matrix transpose. The data access pattern is hence regular. Two programs for blocked LU factorization of a dense matrix form part of the suite. The non-contiguous (LU-n) version has a single producer and multiple consumers. It suffers from considerable fragmentation and false sharing. The contiguous version (LU-c) uses an array of blocks to improve spatial locality. The sorting kernel (RADIX) performs a radix sort on a set of integers. The program suffers from a high-degree of false sharing at page granularity during a permutation phase. Each processor writes the values of its keys into a random location in a shared array. The Ocean simulation program with contiguous partitions (OCEAN-c) simulates large scale ocean movements. This version uses a red-black Gauss-Seidel multi-grid equation solver and has a regular nearest-neighbor type access pattern. The last two programs tested evaluate the forces and potentials occurring over time in a system of water molecules. The first version (WATER-n2) uses a simple data structure, which results in a less efficient algorithm. The second version (WATER-sp) uses a 3-D grid of cells so that a processor that owns a cell only needs to look at neighboring cells to find interacting molecules. Communication arises out of the movement of molecules from one cell to another at every time-step.

4.2 Image Deblurring

The application tested is a parallel algorithm for deblurring of images obtained from Magnetic Resonance Imaging. The images generated may suffer a loss of clarity due to inhomogeneities in the magnetic field. One of the techniques for removing this blurring artifact is the demodulation of the data for each pixel of the image using the value of the magnetic field near that point in space. This method consists of acquiring a local field map, finding the best fit to a linear map and using it to deblur the image distortions due to local frequency variations. This is a very computation intensive operation and has previously been parallelized using a message passing approach [17]. Each thread deblurs the input image around its chosen frequency points in parallel. After deblurring around a frequency point, the relevant portions of the final image have to be updated. Since the portions updated by different threads could overlap, the update has to be done under the protection of a global lock.

4.3 Matrix Multiplication

The matrix multiplication program computes the product of two dense square matrices. The resultant matrix is divided across the processors using a row-wise block dis-

tribution. The size of each of the blocks has been set to a multiple of the page size of the machine. Each application thread computes at least one complete block of contiguous values. Hence there is no false sharing in this implementation and the application is close to ideal for execution on a page based DSM system.

4.4 Successive Over Relaxation

The successive over relaxation program (SOR) is a kernel for solving partial differential equations using a red-black algorithm. In every iteration, each point in a grid is set to the average of its four neighbors. Most of the traffic arises out of nearest neighborhood communication at the borders of the grids.

4.5 Program Parameters

The problem sizes used in this evaluation are shown in Table 2.

Program	Parameters
FFT	1048576 complex doubles
LU-c	2048 × 2048 doubles, block size 128
LU-n	1024 × 1024 doubles, block size 32
RADIX	1048576 integers, radix 1024
OCEAN-c	258 × 258 grid
WATER-n2	4096 molecules, 3 steps
WATER-sp	1000 molecules, 3 steps
MATMULT	1024 × 1024 doubles, 16 blocks
SOR	2002 × 1002 doubles, 100 iterations
MRI	14 frequency points, PHANTOM image

Table 2: Problem Sizes for Network Evaluation

5 Experimental Results

Table 3 shows the utilization of each network in terms of the the bandwidth used, and the average roundtrip latency per node. The average packet size for an application is the ratio of the total bytes transmitted to the number of messages sent. Both these numbers are characteristic of the application itself, and independent of the particular type of network. The bandwidth utilization (BW) was calculated as the ratio of the total bytes transmitted per task to the lifetime of the task. This is a conservative estimate, since some of the applications have a very high computation to communication ratio. The roundtrip latency for a message includes the processing overhead at the remote node.

Programs like LU-c, WATER-n2, WATER-sp, MATMULT, and MRI have low bandwidth requirements and that is reflected in the small performance improvements while going from 10 Mbps Ethernet to the higher bandwidth networks as shown in Figures 3, 4, 5, 6, and 7.

FFT, LU-n, and SOR have a high bandwidth requirement, and there is a distinct network bottleneck when us-

ing 10 Mbps Ethernet. This can be seen in Figures 8, 9, and 12, FastEthernet and 155 Mbps ATM are able to provide a speed-up when more nodes are added to the cluster, while the 10 Mbps cluster shows a slowdown with these programs.

RADIX and OCEAN-c show similar trends in Figures 10 and 11. FFT, LU-c, LU-n, and RADIX have a large average packet size. From Figure 1, the latency of FastEthernet is lower than that of the ATM network only for packets less than approximately 4000 bytes. Hence one can see some performance improvement by using the ATM network in these cases. On the other hand, WATER-n2 and WATER-sp have a very small packet size, in this case the low latency FastEthernet outperforms ATM. For the other programs there is not much difference in the performance of the two networks.

MRI has a very high computation to communication ratio, and the final result is updated under the protection of a lock, hence it shows the least latency in all three cases. MATMULT has similar behavior, but at termination all threads try to update their part of the resultant matrix, which causes the latency to increase due to collisions at the receiving node.

Figure 13 shows the roundtrip times for *Strings* messages on each network, using different payload sizes. These results show that the communication system introduces a fair amount of latency when compared to the raw UDP data. The roundtrip time is measured for request-reply pairs, with no event handling at the server side. Hence these represent upper bounds that the current runtime system may achieve. This indicates that some work needs to be done to improve the performance of the *Strings* communication system

6 Conclusion

Overall the communication results show that for half the applications, namely LU-c, WATER-n2, WATER-sp, MATMULT, and MRI, the bandwidth utilization is not very high, and there does not seem to be much of an advantage from using higher bandwidth networks. However, for applications like FFT, LU-n, and SOR, the 10 Mbps Ethernet network poses a bottleneck. These programs show a speed-up when run on higher bandwidth networks, but a slowdown on 10 Mbps Ethernet.

In the worst case, if all nodes communicate at the same time, on the same link, the total bandwidth requirement is still within the theoretical limits for both high bandwidth networks. However the latencies observed for most programs are much higher than predicted by the raw UDP data. This indicates that there is substantial software overhead in the runtime system.

The study also shows that in the current system, the extra cost of the ATM network and switch rarely translates

Program	Average Msg. Size (bytes)	Ethernet		FastEthernet		OC-3 ATM	
		BW (Mbps)	Latency (ms)	BW (Mbps)	Latency (ms)	BW (Mbps)	Latency (ms)
FFT	4945.50	2.12	134.23	12.71	174.16	15.84	111.68
LU-c	5976.38	1.93	431.08	3.71	192.73	4.68	256.80
LU-n	6223.95	2.02	178.30	9.32	146.84	9.87	148.68
RADIX	4886.12	2.16	32.31	12.85	3.62	14.83	2.80
OCEAN-c	4059.26	1.91	274.97	8.44	36.39	8.42	31.90
WATER-n2	197.19	0.62	118.15	0.77	170.82	0.71	313.30
WATER-sp	893.56	0.92	15.36	1.47	4.69	1.29	5.27
MATMULT	4385.72	1.00	48.57	1.42	7.03	1.45	5.24
SOR	7442.80	2.14	84.69	15.44	327.80	14.41	373.75
MRI	3440.01	2.13	7.31	4.39	1.70	4.40	1.62

Table 3: Communication Utilization Per Node (16 tasks)

into a performance improvement compared to the FastEthernet. Out future work would include studying the behavior of the programs when using switched FastEthernet, as well as switched Gigabit Ethernet.

Finally, the *Strings* communication subsystem introduces a large latency when compared to raw UDP packets. We plan to isolate these overheads more precisely, and investigate better approaches to providing a fast reliable communication system for the DSM.

References

- [1] C. Amza, A. L. Cox, S. Dwarkadas, P. Keleher, H. Lu, R. Rajamony, W. Yu, and W. Zwaenepoel, "TreadMarks: Shared Memory Computing on Networks of Workstations," *IEEE Computer*, pp. 18–28, February 1996.
- [2] D. Khandekar, *Quarks: Portable Distributed Shared Memory on Unix*. Computer Systems Laboratory, University of Utah, beta ed., 1995.
- [3] P. Keleher, *CVM: The Coherent Virtual Machine*. University of Maryland, CVM Version 2.0 ed., July 1997.
- [4] R. Stets, S. Dwarkadas, N. Hardavellas, G. Hunt, L. Kontothanassis, S. Parthasarathy, and M. Scott, "CASHMERE-2L: Software Coherent Shared Memory on a Clustered Remote-Write Network," in *Proceedings of the ACM Symposium on Operating System Principles*, (Saint Manlo, France), October 1997.
- [5] S. Roy and V. Chaudhary, "Strings: A High-Performance Distributed Shared Memory for Symmetrical Multiprocessor Clusters," in *Proceedings of the Seventh IEEE International Symposium on High Performance Distributed Computing*, (Chicago, IL), pp. 90–97, July 1998.
- [6] S. C. Woo, M. Ohara, E. Torri, J. P. Singh, and A. Gupta, "The SPLASH-2 Programs: Characterization and Methodological Considerations," in *Proceedings of the International Symposium on Computer Architecture*, pp. 24–36, June 1995.
- [7] J. B. Carter, "Design of the Munin Distributed Shared Memory System," *Journal of Parallel and Distributed Computing*, 1995.
- [8] R. Mirchandaney, S. Hiranandani, and A. Sethi, "Improving the Performance of DSM Systems via Compiler Involvement," in *Proceedings of Supercomputing 1994*, 1994.
- [9] Information Networks Division, *Netperf: A Network Performance Benchmark*. Hewlett-Packard Company, <http://www.netperf.org>, revision 2.1 ed., February 1996.
- [10] A. L. Cox, S. Dwarkadas, P. Keleher, H. Lu, R. Rajamony, and W. Zwaenepoel, "Software vs. Hardware Shared Memory Implementation: A Case Study," in *Proceedings of the International Symposium on Computer Architecture*, pp. 106–117, April 1994.
- [11] C. Amza, A. Cox, K. Rajamani, and W. Zwaenepoel, "Tradeoffs Between False Sharing and Aggregation in Software Distributed Shared Memory," in *Proceedings of the ACM Symposium on the Principles and Practice of Parallel Programming*, (Las Vegas), pp. 90–99, June 1997.
- [12] H. Lu, A. L. Cox, S. Dwarkadas, R. Rajamony, and W. Zwaenepol, "Compiler and Software Distributed

Shared Memory Support for Irregular Application,” in *Proceedings of the ACM Symposium on the Principles and Practice of Parallel Programming*, 1997.

- [13] P. Keleher and C.-W. Tseng, “Enhancing Software DSM for Compiler-Parallelized Applications,” in *Proceedings of International Parallel Processing Symposium*, August 1997.
- [14] R. Hempel, H. C. Hoppe, U. Keller, and W. Krotz, *PARMACS V6.0 Specification. Interface Description*. Pallas GmbH, November 1993.
- [15] D. Jiang, H. Shan, and J. P. Singh, “Application Restructuring and Performance Portability on Shared Virtual Memory and Hardware-Coherent Multiprocessors,” in *Proceedings of the ACM Symposium on the Principles and Practice of Parallel Programming*, (Las Vegas), pp. 217 – 229, ACM, 1997.
- [16] Y. Zhou, L. Iftode, J. P. Singh, K. Li, B. R. Toonen, I. Schoinas, M. D. Hill, and D. A. Wood, “Relaxed Consistency and Coherence Granularity in DSM Systems: A Performance Evaluation,” in *Proceedings of the ACM Symposium on the Principles and Practice of Parallel Programming*, (Las Vegas), pp. 193 – 205, June 1997.
- [17] P. Menon, V. Chaudhary, and J. G. Pipe, “Parallel Algorithms for deblurring MR images,” in *Proceedings of ISCA 13th International Conference on Computers and Their Applications*, March 1998.

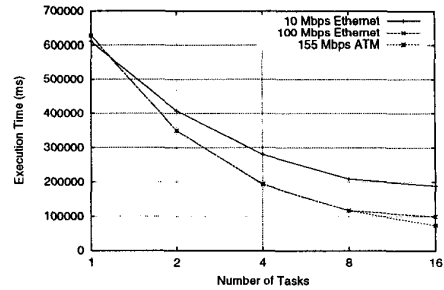


Figure 3: LU-c

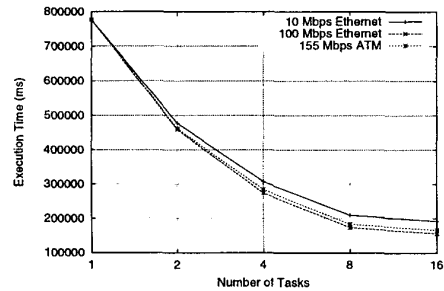


Figure 4: WATER-n2

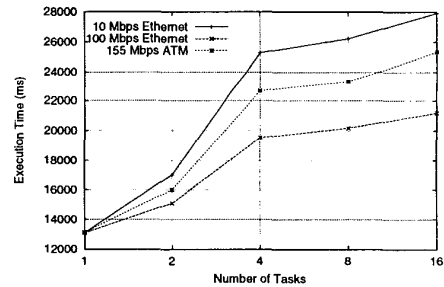


Figure 5: WATER-sp

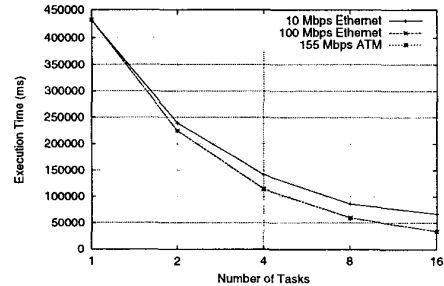


Figure 6: MATMULT

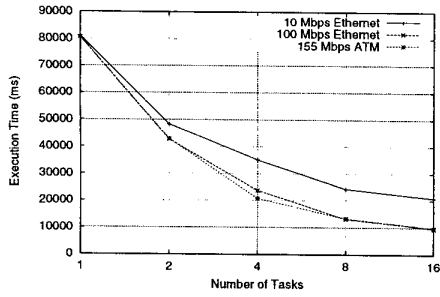


Figure 7: MRI

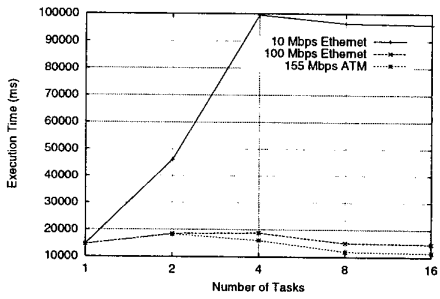


Figure 8: FFT

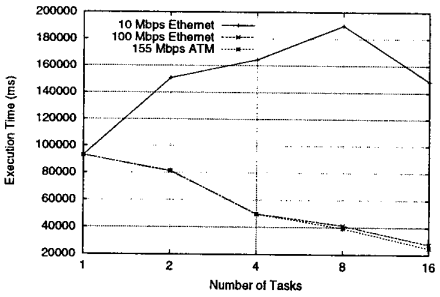


Figure 9: LU-n

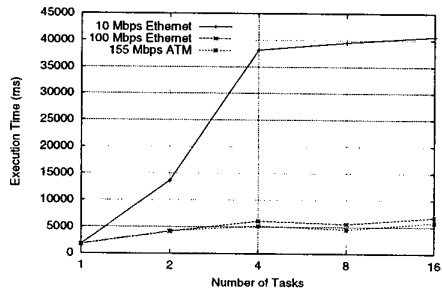


Figure 10: RADIX

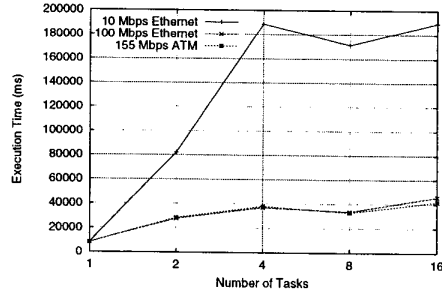


Figure 11: OCEAN-c

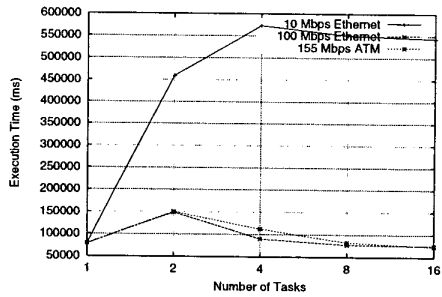


Figure 12: SOR

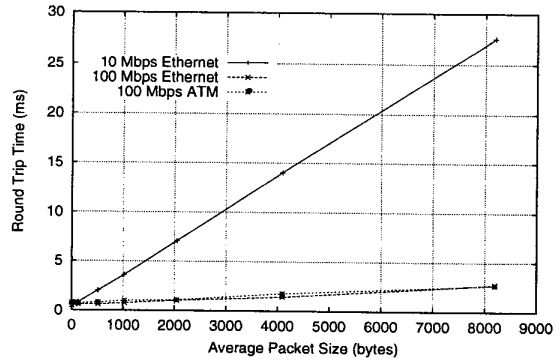


Figure 13: Roundtrip times for *Strings* messages