



## Using Warp to Control Network Contention in Mermera\*

Abdelsalam Heddaya  
heddaya@cs.bu.edu

Kihong Park<sup>†</sup>  
park@cs.bu.edu

Himanshu Sinha<sup>‡</sup>  
hss@cs.bu.edu

BU-CS-93-007

June 4, 1993

Revised September 25, 1993

To appear in Proc. 27th Hawaii Int'l Conf. on System  
Sciences, IEEE Comp. Soc. Press, January 1994.

### Abstract

Parallel computing on a network of workstations can saturate the communication network, leading to excessive message delays and consequently poor application performance. We examine empirically the consequences of integrating a flow control protocol, called Warp control [Par93], into Mermera, a software shared memory system that supports parallel computing on distributed systems [HS93].

For an asynchronous iterative program that solves a system of linear equations, our measurements show that Warp succeeds in stabilizing the network's behavior even under high levels of contention. As a result, the application achieves a higher effective communication throughput, and a reduced completion time. In some cases, however, Warp control does not achieve the performance attainable by fixed size buffering when using a statically optimal buffer size.

Our use of Warp to regulate the allocation of network bandwidth emphasizes the possibility for integrating it with the allocation of other resources, such as CPU cycles and disk bandwidth, so as to optimize overall system throughput, and enable fully-shared execution of parallel programs.

**Keywords:** Distributed non-coherent shared memory, network contention, flow control, iterative methods, Isis.

---

\*This research was supported in part by NSF under grants IRI-8910195, IRI-9041581 and CDA-8920936.

<sup>†</sup>This author was additionally supported by NSF grant CCR-9204284.

<sup>‡</sup>Present address: GTE Laboratories, 40 Sylvan Road, Waltham MA 02254.

# 1 Introduction

Today's high performance networked workstations are increasingly used for parallel computing. Such systems are connected by networks (*e.g.*, Ethernet, FDDI) that are proportionately slower than parallel computer interconnection networks, in terms of latency and bandwidth. Consequently, parallel programs can place very high sustained communication loads on workstation networks, causing congestion and high communication delays. This phenomenon severely limits the number of workstations that can participate fruitfully in parallel when running communication-intensive applications.

Mermera is an example of a system that encourages applications to generate high traffic volumes, causing network congestion. It is a software shared memory system that has been developed to support parallel computing on a workstation network [HS93, Sin93]. Processes that comprise a parallel program communicate through read and write operations to the shared memory provided by Mermera. Several different memory behaviors are supported, including one coherent<sup>1</sup>, and three non-coherent (Pipelined Random Access Memory [LS88], Slow Memory [HA90], and Local Consistency [HS93]).

A key property of non-coherent memory operations supported by Mermera is that they are “asynchronous” in that they do not wait for any communication to be performed before they return. This allows the application program to issue as many shared memory operations as the processing power of the CPU allows, potentially flooding the network with communication requests. By contrast, coherent write operations are inherently “synchronous” in that they do not return until some message is successfully sent and received. Applications such as asynchronous iterative algorithms [BT89] and oblivious computations [LS88], are thus capable of causing high and sustained network contention. Operating systems that aim to support parallel computing, must therefore offer mechanisms for controlling network contention to a higher degree than is available today. To study the question, and to try to formulate requirements for the likely candidate mechanisms for this purpose, we implemented, tested and measured the Warp control protocol under Mermera.

*Warp control* [Par93] is a distributed end-to-end flow control protocol which uses a time-stamp based estimator of network state to adaptively throttle arrival rates so as to optimize network throughput. It consists of a protocol suite, of which *rate adjustment protocol* (RAP) is the basic one used in the experiments. The network estimator is called *warp*, and it uses time-stamps to

---

<sup>1</sup>Our definition of coherence, first advanced in [HS92] is equivalent to that of sequential consistency [Lam79].

estimate relative change in network delay. We will often use “Warp” (capitalized) to refer to the protocol itself. In case of single server queues,  $warp \approx \rho$  where  $\rho$  is the network utilization (*i.e.*, throughput). Warp control has the desirable property of being asymptotically stable, although for networks characterized by unimodal load–service rate functions<sup>2</sup>, additional controls are needed to ensure stability. Warp is particularly suited for steady, high-bandwidth traffic, a communication pattern supported by non-coherent writes in Mermera. Although the ideal place for Warp to reside is in the network interface unit (*e.g.*, network layer in the ISO-OSI hierarchy or lower), in the current implementation it was applied at the level of Mermera which itself runs on top of Isis. Due to Isis’ own flow control algorithm and Ethernet’s exponential back-off mechanism, Warp was receiving highly filtered information, but nevertheless achieving improved performance when running a communication-intensive application. In order to harness the computational resources of large networks of workstations effectively, we believe that the network must be treated as an integral part of the distributed system, to be managed in tandem with other resources. We elaborate on this point in sections 3 and 7.

In the next two sections, we describe the essential details of Mermera and Warp. Section 4 specifies our experimental set-up, and section 5 provides basic measurements that characterize the communication underlying Mermera with, and without, Warp. Preliminary measurements that exemplify the application performance improvements obtainable with Warp are given and interpreted in section 6. Section 7 discusses our conclusions.

## 2 Mermera

Mermera supports several different memory behaviors, including:

- **Coherent Memory:** All processes agree on the order of *all* writes, a definition that can be shown to be equivalent to sequential consistency [Lam79]. More precisely, if a process observes some writes in a certain order then no other process observes those writes in a different order. A process  $i$  is said to observe a write  $x.w$  to location  $x$ , if the value returned to process  $i$  by a read operation is the value written by  $x.w$ , or is a value written by a process that observed  $x.w$ .

---

<sup>2</sup>One method for modeling network congestion is to assume a functional relationship between network load and service rate that has a unimodal shape [Par93].

- **Pipelined RAM (PRAM):** This behavior shares only the acronym with the Parallel Random Access Machine, a theoretical model of parallelism often used to design simple parallel algorithms. PRAM requires that the order of all writes by the *same* process is respected by all processes, *i.e.*, if a process performs two writes,  $w_1$  followed by  $w_2$ , then no process can read them in the reverse order. Writes by different processes may be interleaved in different orders by different processes.
- **Slow Memory:** All writes by *the same process to the same location* are ordered by all processes in the order they were written. Writes to different locations by the same process may be ordered differently by different processes.

Our current implementations<sup>3</sup> support the above behaviors using full replication of the shared memory. The interface consists of one read operation, and several write operations, one for each type of memory behavior. A read operation returns the value stored in the local copy of shared memory, and hence takes the same time as an ordinary memory reference. This time is typically 0.1–1.0 microseconds. Each write operation causes Mermera to enqueue the new value for broadcast over the network. A coherent write operation returns only after the message has been sent and received back by the writer.

Slow and Pipelined RAM write operations do not wait for any communication to complete; once their request for communication is recorded locally, the operation can return. These updates can be enqueued for later transmission when an appropriate number of them can be batched together for efficiency. We can say that these operations are asynchronous with respect to communication, and computations that use them can, in principle, submit new values for transmission at a sustained rate greater than the network capacity.

Writes are propagated by broadcasting the values to other processes. We use the Isis toolkit [BJ87, BSS91] for our implementation because it provides a suite of multicast protocols that satisfy different ordering properties. The broadcasts of interest to us are *abcast()*, *fbcast()* and *mbcast()*. These primitives have different constraints on the order in which the messages are delivered to their destinations.

All messages sent using *abcast()* are delivered in the same order at all destinations, *i.e.*, the order in which these messages are delivered is the same for all processes. This is exactly the property we want for *CO\_Write*. The *fbcast()* messages obey a weaker constraint: messages sent by the

---

<sup>3</sup>Mermera has been implemented on a BBN Butterfly, on a CM-5, and on a workstation network. Congestion control is applicable only in the latter.

*same* process are delivered to all processes in the order they were sent. However, *fbcasts* sent by different processes may be interleaved in different orders at different recipients. This suffices for *PRAM\_Writes*, so we use *fbcasts* to propagate them. No ordering constraints are guaranteed among *mbcast()* messages. We use this primitive for propagating *Slow\_Writes*.

Isis offers no ordering guarantees among messages sent using different primitives, *e.g.*, if two messages are sent one after the other using *abcast()* and *fbcast()*, respectively, they are not necessarily delivered in the order they were sent. Our implementation enforces this ordering explicitly via sequence numbers, and by careful control of the order in which received updates are applied.

Isis employs the UDP transport protocol, which, in turn, runs on top of IP and Ethernet's CSMA/CD protocol. Both Isis and the Ethernet, but not UDP, have built-in flow control whose effect can be discerned in our measurements presented in sections 5 and 6.

Numerous software shared memory systems have been proposed in the literature, but most of them provide coherence at the level of sequential consistency or stricter. Recently, some systems started to allow weaker conditions. Examples include the DASH multiprocessor [LLJ<sup>+</sup>93], the Munin software shared memory system [CBZ91], and transactional memory [HM92], which allow sequential consistency to be broken, but only during the execution of code blocks explicitly marked by the programmer. These systems, as well as Mermera, stand to benefit from network contention control protocols.

### 3 Warp Control

Warp control is a distributed end-to-end flow control protocol that uses a time-stamp based scheme to throttle arrival rates at active message sources to achieve optimal network utilization. For a related discussion of flow control algorithms, see [GK80, HW91]. Let  $N$  be the network, and let  $p_1, p_2, \dots, p_n$  be the network nodes. Let  $\mu$  be the service rate of the network, and let  $\lambda_1, \lambda_2, \dots, \lambda_n$  be the arrival rates of  $p_1, p_2, \dots, p_n$ , respectively. Thus the total arrival rate of the network is given by  $\lambda = \sum_{i=1}^n \lambda_i$ , and the utilization can be defined as  $\rho = \lambda/\mu$ . For simplicity of exposition, assume that every node  $p_i$  has access to the same global clock<sup>4</sup>. Assume at time  $t_1$ , node  $p_i$  wants to send a message  $m_{ij}$  to node  $p_j$ . The following protocol is executed at the sender  $p_i$ :

#### Message encoding protocol (MEP):

---

<sup>4</sup>In general, this assumption can be relaxed so that, not only do the network nodes need not have access to the same absolute time, but their local clocks may also run at different speeds.

1. Create a header containing  $i, j$ , and time stamp  $t_1$ , that is,  $m_{ij}.time\_stamp := t_1$ .
2. Attach the header to the data section of  $m_{ij}$  and submit to network  $N$ .

For all  $k \in \{1, 2, \dots, n\} \setminus \{j\}$ ,  $p_j$  maintains a data structure  $hist[k]$  containing two fields  $hist[k].last\_in$  and  $hist[k].last\_out$ . In  $hist[k].last\_in$  is recorded the time at which the last message from  $p_k$  arrived at  $p_j$ , and  $hist[k].last\_out$  records the time at which it was sent out from  $p_k$ . When  $p_j$  receives  $m_{ij}$ , it executes the following protocol, in which  $t_2$  denotes the message arrival time.

**Message decoding protocol (MDP):**

1.  $warp := \frac{t_2 - hist[i].last\_in}{m_{ij}.time\_stamp - hist[i].last\_out}$
2.  $hist[i].last\_in := t_2$
3.  $hist[i].last\_out := m_{ij}.time\_stamp$

It can be shown [Par93], under the assumption of  $N$  being a FCFS single-server queue and  $|d\lambda/dt| \ll 1$ , that  $warp \approx \rho$ . Therefore, to achieve and maintain maximum utilization, the following control law can be employed,

$$\frac{d\lambda}{dt} = \epsilon(1 - warp) \tag{1}$$

where  $\epsilon > 0$  is a constant adjustment factor. If  $\tau$  is the network delay, one can further show that the resulting system is asymptotically stable if

$$0 < \tau \frac{\epsilon}{\bar{\mu}} < \frac{\pi}{2} \tag{2}$$

where  $\bar{\mu}$  is the mean service rate. Thus to attain maximum utilization in a stable manner, it suffices to set  $\epsilon$  very small. In practical terms, this translates to an initial transient period at start-up time, and a lagged response to sudden structural changes (*e.g.*, 30% of the nodes going down at once). On the positive side, this makes the system robust with respect to noise, guaranteeing asymptotic stability even under dynamic perturbation. The issues of better responsiveness, fairness and stability under unimodal load-service rate functions are discussed in [Par93].

Distributed control is achieved by noting that  $\lambda = \sum_{i=1}^n \lambda_i$ . That is, at every  $p_j$ , the following update protocol is executed after performing MDP:

**Rate adjustment protocol (RAP):**

1.  $\lambda_j(t) := \lambda_j(t - 1) + \epsilon_j(1 - warp)$ .

2. If  $\lambda_j(t) < 0$ , then  $\lambda_j(t) := 0$ .

A straightforward way of implementing message submission rate control at every  $p_i$  using a *buffering scheme* is to maintain a large message queue  $B$ , and enqueue message submissions in  $B$ . The network interface unit at  $p_i$  periodically examines  $B$ , and at every such instant dequeues  $\lambda_i b$  messages from  $B$ . It executes MEP on the dequeued messages and sends them off. The parameter  $b > 0$  is system dependent, and affects the quantization range of  $\lambda_i$ . Information regarding the current state of  $B$  and  $\lambda_i b$  can be made available to the application running at  $p_i$ , so that it can better manage its communication needs should its computation allow it to do so. For example, in asynchronous iterative algorithms that solve fixed-point problems arising in various applications [BT89], convergence rate is enhanced if the application can keep its computation/communication ratio, which is variable, compatible with  $\lambda_i b$ . This way, no time is wasted in generating messages that may get delayed due to unnecessary buffering. For applications that do not have controllable computation/communication ratios, this benefit cannot be harnessed.

Nevertheless, from an operating systems point of view, if the operating system is managing a set of processes each with its own nonvariable computation/communication ratio, the operating system itself may schedule the processes so that the system as a whole stays as compatible as possible with  $\lambda_i b$ . For example, if  $\lambda_i b$  is small, then it is wasteful for the operating system to allocate CPU cycles to processes having a small computation/communication ratio (*i.e.*, communication intensive) if it has a choice. Not only will the delay be large due to the growing queue, but in the extreme case when  $B$  is full, CPU cycles will be wasted trying to write to a buffer that is already full. These speculative scenarios illustrate the close relationship between various resources such as computation (CPU) and communication (network) in networked workstations that share them. Thus, the possibility clearly exists for an operating system to enhance system performance by allocating CPU and network resources in tandem.

In the current implementation, due to issues specific to the present design of Mermera, a simplified scheme is being used. A buffer  $B$  is maintained between Mermera and Isis, and if  $|B| > \Lambda_i b$ , where  $|B|$  denotes the queue length, the buffer is flushed to Isis.  $\Lambda_i$  is a quantity that has the opposite effect of  $\lambda_i$ , in the sense that, if  $\Lambda_i$  increases, the message submission rate decreases due to less frequent flushing. This modified implementation is easily accommodated by changing the sign of  $\epsilon$  in equation (1) from positive to negative. The value referred to as “lambda” in this paper is actually  $\Lambda$ , and should not be confused with the arrival rate  $\lambda$ .

One can view  $\Lambda$  as a parameter that controls the amount of bandwidth allocated to the com-

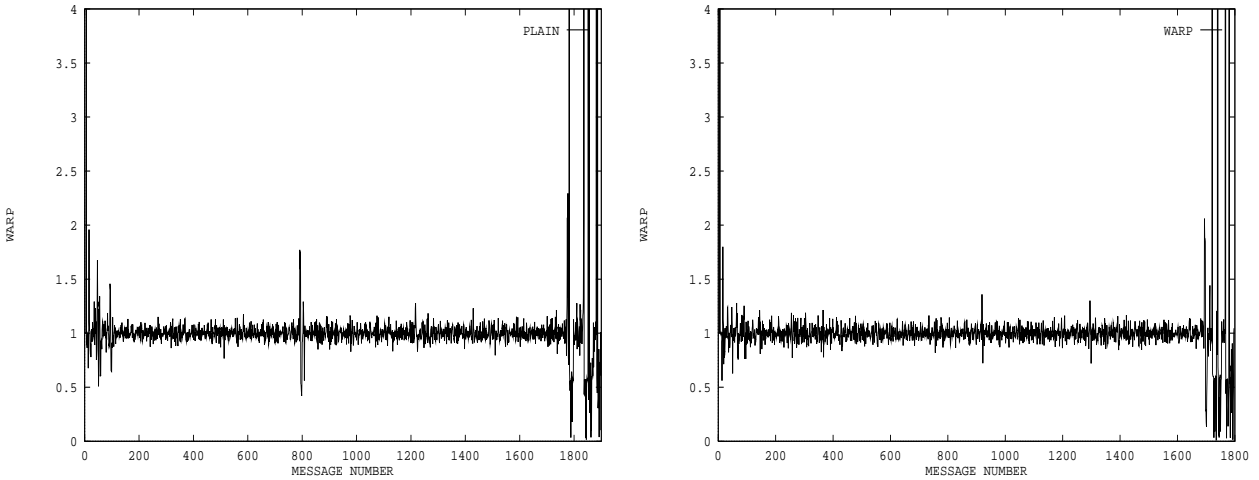


Figure 1: In the absence of background traffic, *warp* values indicate that the network is well behaved both without Warp (left) or with it (right). The explosion in delay represented by the high warp peaks toward the end of the computation signify the extra communication our example application needs to detect termination.

putation whose messages are “throttled” according to the value of  $\Lambda$ . Under this interpretation, mechanisms that tie CPU scheduling to network bandwidth allocation can be exemplified by the following scheme. Whenever a process attempts to send a message, and finds that its message will have to be queued for a long time because  $|B| \gg \Lambda b$ , the system gives the process the following choice. It can either give up the CPU in favor of other processes that may have some spare communication bandwidth allocated to them, or it can abort the sending of its message, hoping to generate another message with more current or complete information by the time its share of bandwidth permits the transmission of the message. This scheme has meaning under a more general Warp Control system than is implemented in this paper, one that allocates bandwidth among different parallel programs that share the system’s communication facility.

## 4 Experimental Set-up

We conducted our measurements on a network of six dedicated Sun Sparc 1+ workstations and a server, each equipped with 660 MByte individual disks, and connected via a private 10 Mbit Ethernet. The workstations all run SunOS version 4.1.1, and the Isis Toolkit version 2.2.5. We ran NetMetrix version 3.0.1 on the server and one of the workstations, to generate background traffic

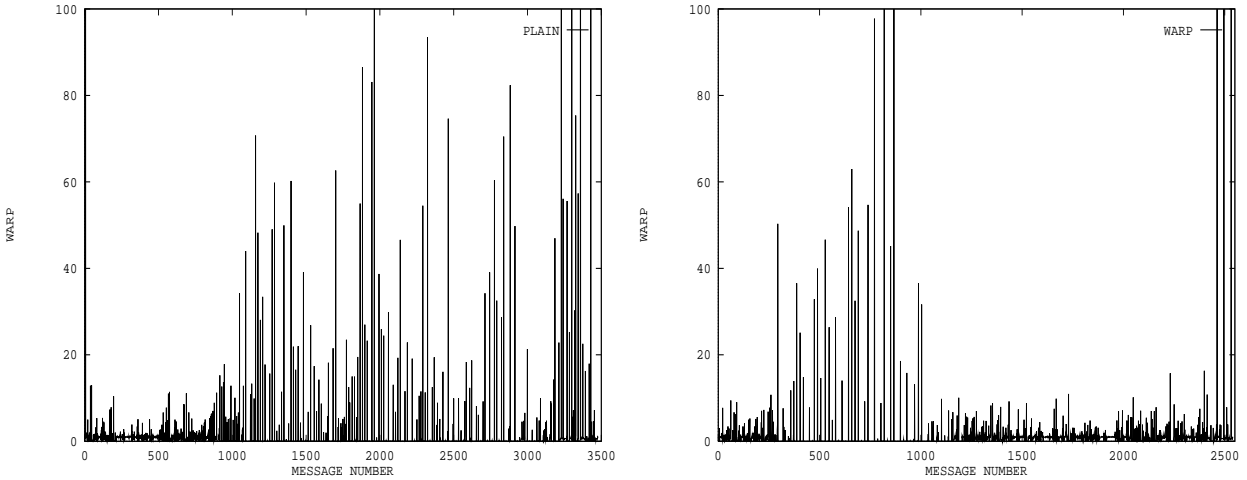


Figure 2: At 65% background traffic, *warp* values exhibit numerous high spikes (above 10) without Warp control (left). Warp control (right) effectively eliminates them, once it stabilizes. Spikes above 100 are clipped to preserve scale.

in order to increase contention for the Ethernet.<sup>5</sup> The remaining five workstations were dedicated to running the parallel equation solver on top of Mermera and Isis.

We measured the completion time of the equation solver on a system of linear equations  $Ax + b = 0$  in 1000 variables, generated pseudo-randomly. The level of traffic infused by NetMetrix on each of the two machines, given as an uninterpreted percentage, was varied between 0% to 65%. We also measured the input and output quantities for Warp control, *warp* and  $\Lambda$ , defined in the previous section.

Warp was implemented as a function call from Mermera, invoked each time a message is received. The corresponding message rate control is incorporated as a conditional check, tested at every message enqueue, flushing the buffer to Isis if  $|B| > \Lambda_i b$ . This places Warp control at an inefficiently high level, acting on data that is filtered through lower layer protocols. For operating system support purposes, Warp should be moved to the transport/network layer of the ISO-OSI reference model.

## 5 Network Performance

Figure 1 shows the variations in network delay (as measured by *warp*) during a typical run of the equation solver on top of Mermera. With only five workstations, the network is clearly operating

---

<sup>5</sup>We are grateful to Metrix Network Systems, Inc., for free use of their software.

well within its capacity, and therefore is very well behaved.<sup>6</sup> Turning on Warp control does not lead to any noticeable performance improvement. In order to see what might happen if the network contention was high, for example, caused by having 100 workstations instead of only five, we used a traffic generation program that is part of the NetMetrix software. We ran the traffic generator on two dedicated workstations, other than the five that run the equation solver, and obtained the results shown in figure 2.

Compared to figure 1, *warp* values in figure 2 are far higher, with many spikes exceeding 5, some exceeding 50, and a few exceeding 500. A value of 5 means that a Mermera internal message took five times longer to reach its destination than the immediately preceding message from the same source to the same destination. Since Mermera messages, which are sent via Isis, take roughly 10-50 ms on average to reach their destination, our measurements suggest that some messages are taking 0.5-2.5 seconds—an extremely long delay. Such large differences in message delay can be explained by looking at the flow control mechanisms used by Isis and the underlying Ethernet. Isis uses several schemes, among them a simple heuristic that, based on the number of buffered messages at a node, determines whether congestion exists or not. Once Isis decides that congestion exists, it stops sending messages for a long period of time (1 second) [Bir93], thus causing a large spike in *warp*. Smaller spikes may be explained by the Ethernet’s back-off mechanism when collisions occur [Tan88].

Figure 2 offers empirical insight into the behavior of Warp: it slowly but effectively squelches large spikes in *warp*, and consequently in message delay. We see in the right hand plot that the spikes subside after a thousand messages or so have been delivered. Warp’s success in smoothing the network’s behavior reduces the total number of messages needed to run the application to completion, from  $\sim 3500$  to  $\sim 2500$ . Both of these observations are further elaborated in the next section.

## 6 Application Performance

### 6.1 Six-Workstation Experiments

Our preliminary measurements suggest that Warp control can improve application performance by enabling it to extract a higher communication throughput, with a correspondingly smaller delay,

---

<sup>6</sup>The very large spikes at time 1800 are caused by a barrier mechanism that checks for a global termination condition.

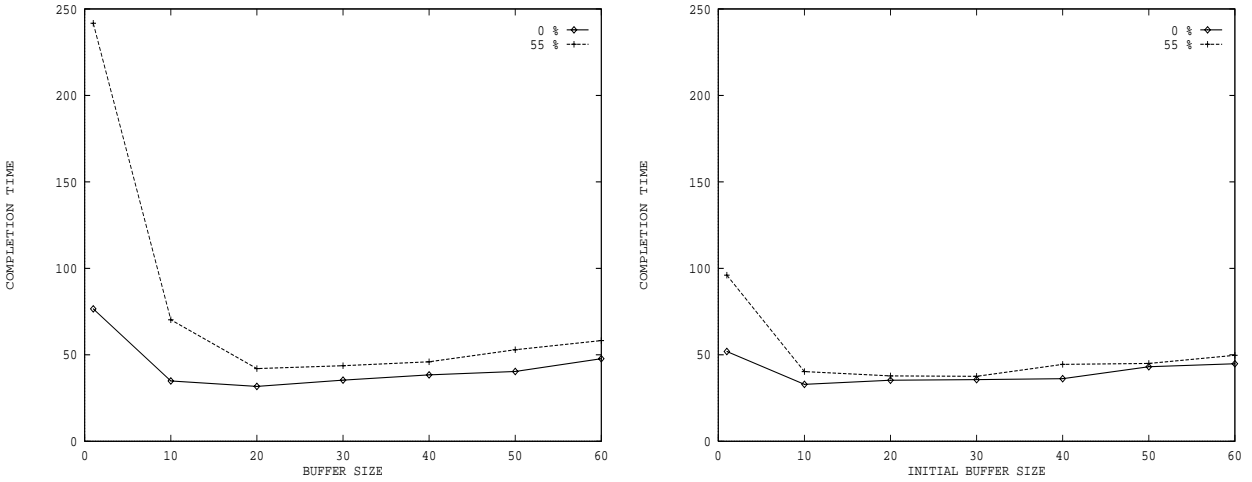


Figure 3: Warp (right) reduces the sensitivity of completion time to background traffic level, compared to the case without Warp (left). The solid, lower, curve corresponds to no background traffic (0%).

from the network. In this section we present evidence to support this conclusion, and empirically analyze the behavior of Warp control in order to explain the precise nature of the performance improvement it affords. In particular, we seek to illuminate the relationship between performance under statically optimal buffer size, and that under dynamically optimal buffer size, controlled by the Warp protocol.

Consider the two different views of the completion time versus the buffer size (initial buffer size when Warp is active), found in figures 3 and 4. Figure 3-left illustrates the detrimental effect network contention can have on the equation solver’s completion time, which is manifested as a shift in the completion time curve. When there is high background traffic and Warp is turned off, simply increasing the fixed buffer size causes a reduction in network contention that is sufficient to decrease the completion time drastically (from about 250 seconds to less than 50 seconds). The vertical separation between the two curves in the same plot (figure 3-left) indicates the impact of background traffic on completion time, which can also be very high when the solver is itself generating high traffic (small buffer size). Warp, when activated, reduces the sensitivity of completion time to the amount of background traffic that exists on the network, suggesting that Warp helps the equation solver extract sufficient throughput from the network even in the presence of 55% background traffic. That is, Warp is able to insulate the system, to a significant degree, from the otherwise harmful effects of increased traffic.

Figure 4 allows us to quantify the performance gains obtained with Warp. The graphs are self-

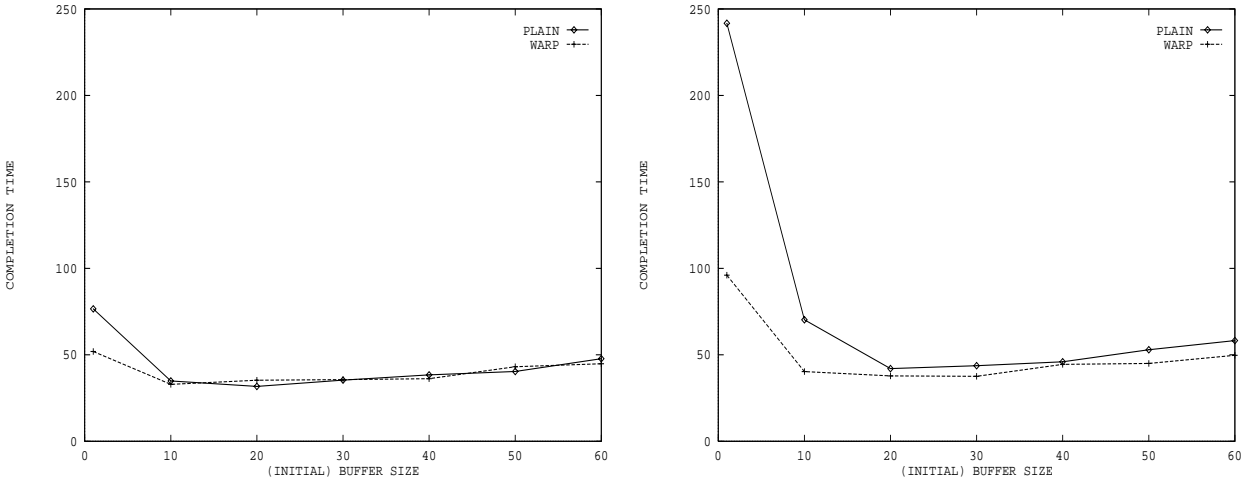


Figure 4: Completion time with five compute processors and two traffic generation processors with and without Warp control. Left and right graphs reflect 0% and 55% background traffic levels, respectively.

explanatory, showing that Warp can relieve the application programmer, or Mermera’s designer, from the need to choose a buffer size appropriate to the particular level of contention that happens to be in effect on the network. Furthermore, in an ever changing environment, Warp provides a simple means for effective adaptation. This suggests the need for operating system tools supportive of such a mechanism, to insulate the programmer from system-dependent communication effects.

However, a very revealing measurement is not apparent in these graphs. Consider the case for an (initial) buffer size of 10, and background traffic level of 65% (only slightly different from the 55% case shown in figure 4-right). The no-Warp run exchanges 34,436 messages, and completes in 95.2 seconds, for a communication throughput of **3,625** updates/second. By contrast, the run under Warp transmits 27,435 messages, and finishes in 52.6 seconds, achieving the *much higher throughput* of **5,276** updates/second. In effect, by reducing the variation in network delay, via controlling excessive traffic, effective throughput is increased, which leads to faster overall performance.

One obvious remaining question is: why does Warp control fail to achieve the best completion time, when it starts with an initial buffer size that is far from optimal? To answer this question, and offer additional insight into the network dynamics under Warp, we present the results in figure 5, which show the actual buffer sizes computed by Warp control throughout two executions. The left hand side pair of graphs correspond to an initial buffer size of 1, while the right hand side pair depict the case for initial buffer size of 10. Clearly, Warp control stabilizes much more slowly when  $\Lambda$  starts farther from its eventual operating neighborhood. In the first case, it takes about 3000 messages for

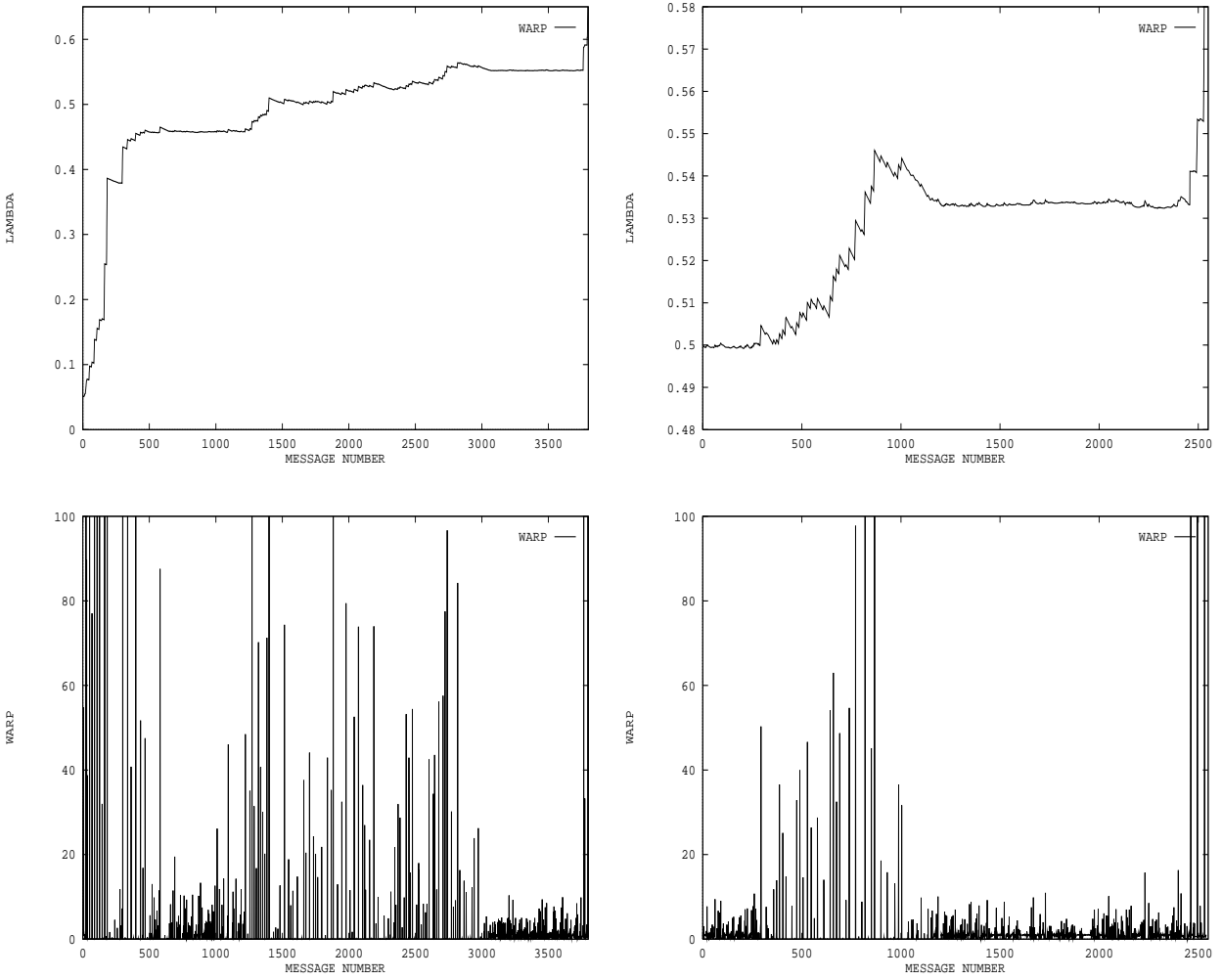


Figure 5:  $\Lambda$  (upper graphs) and *warp* (lower graphs) for background traffic level of 55% with initial  $\Lambda = 0.05$  (left graphs); and traffic 65% with initial  $\Lambda = 0.5$  (right graphs). Buffer size used for each message sent is  $20\Lambda$ , hence the smallest effective change in  $\Lambda$  is 0.05.

Warp to stabilize, while in the second case, it does so after only 1000 messages. The general slowness with which Warp reacts has the advantage of ensuring its long term stability, but, for applications running for short durations, the initial transient period may not get sufficiently amortized to yield significant performance enhancement. On the other hand, for very large problem sizes or when Warp is under the direct control of the operating system, the transient effect will be amortized across multiple applications, yielding an even more pronounced difference in performance.

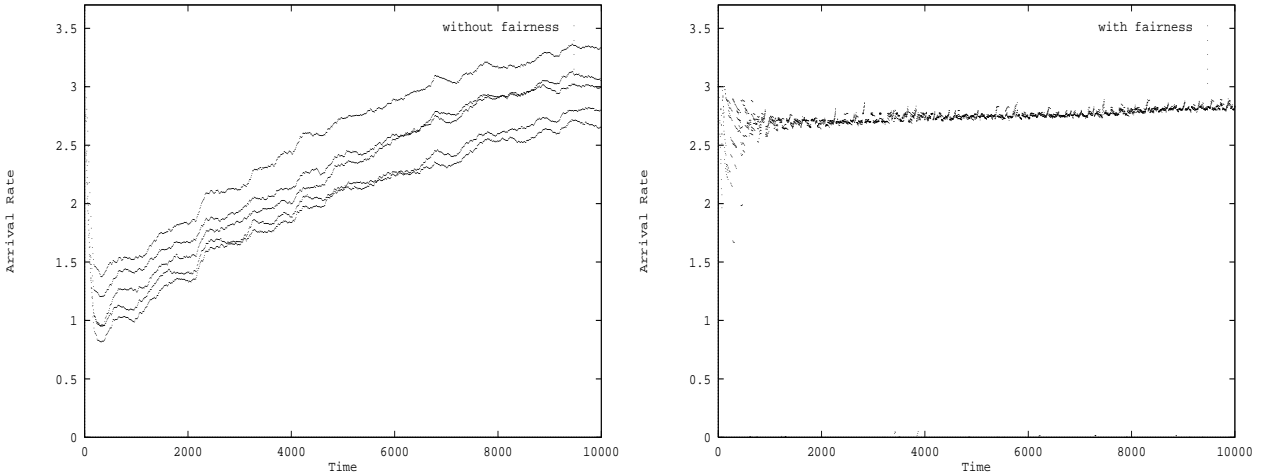


Figure 6: Simulation of a 10 node system under Warp without (left), and with (right), fairness protocol added.

## 6.2 Large-scale Predictions

A relevant phenomenon that is specific to our particular environment concerns the general flatness of the completion time versus buffer size curve, shown in figure 3. The flatness of the curve above a buffer size of 20 indicates the noncritical nature of selecting a near optimal buffer size for the present system, the only requirement being that it not be too small. The flatness is caused in part by the private Ethernet (no routing) and its CSMA/CD collision control mechanism which favors long packets for achieving high throughput. When the number of workstations participating in the parallel computation is large, then the statically optimal buffer size may be much more narrowly defined, making the selection of the optimal buffer size a more challenging task.

A principal reason for this lies with the issue of fairness of bandwidth allocation among participating processors. The completion time of an application distributed over  $n$  nodes depends on the completion time of the slowest component. Hence to ensure timely and sufficient progress over all nodes, bandwidth fairness needs to be imposed, which in turn can be effected by minimizing the variation of  $\Lambda_i$ ,  $E_i\{(\Lambda_i - \bar{\Lambda})^2\}$ , in time where  $\bar{\Lambda} = \sum_i \Lambda_i/n$ . Figure 6 shows a simulation run of a 10 node system without a fairness protocol (left figure), and with one (right figure). The output is taken from [Par93] and shows the variation in bandwidth allocation when fairness is not regulated. The message generation rates of the first 5 nodes are shown, and the drift in individual arrival rates is clearly visible. For a discussion of the fairness protocol, see [Par93].

A prerequisite for reducing variation and maintaining even progress in computation over all

participating nodes is that packets sent out from the nodes not be too long. This implies a more sensitive dependence of completion time on buffer size which increases the need for a dynamic control mechanism such as Warp.

Furthermore, one can define communication uniformity as the level of uniformity of the communication pattern among processors; for example, the higher the locality of communication, the lower the uniformity. Communication uniformity, which grows quadratically in the number of nodes for applications such as asynchronous fix-point iterative algorithms, necessitates shorter packets leading to a more convex-shaped completion time vs. buffer size curve. Due to these aggregate effects, we believe that the role of network bandwidth management for implementing parallel computations on distributed systems will become more critical as the system is scaled up.

## 7 Discussion

This paper deals with a central phenomenon that was brought to our attention through experience with parallel applications running on top of Mermera, namely, extreme changes in message delay and degradation in communication throughput as a result of network contention. Clearly, this phenomenon should not be left to application programmers, or even presentation layer developers, to contend with. For example, neither the equation solver programmer, nor the Mermera designer, have the necessary information to be able to dynamically adjust their software's communication requirements to the levels that would optimize network performance in a changing shared environment. Unfortunately, in the absence of specific operating system support for a flow control protocol such as Warp at lower layers of the ISO-OSI hierarchy, these programmers are left with no other alternative but painstakingly tune their applications by hand, *and* live with performance that can be optimal only under very restricted circumstances, if at all.

We have shown that in systems like Mermera, with a pronounced need for flow control as a result of the “asynchrony” of their communication, that Warp can help smooth network behavior, and deliver enhanced performance to applications. Our measurements also give experimental insight into the dynamic behavior of Warp itself, confirming the theoretical predictions and large-scale simulations of [Par93], given for a general packet or circuit-switched queueing system. Furthermore, our data revealed to us an undocumented aspect of the Isis toolkit, which is the technique it uses to control congestion. It was instructive for us to note that in every run, Warp was able eventually to smooth network traffic enough to prevent Isis' own mechanism from kicking in.

One can think of Warp control as a bandwidth allocation mechanism. This raises many questions, briefly touched upon in section 3, including the possibility of tying CPU scheduling to network bandwidth allocation in order to improve total system throughput. Such a “holistic” perspective has obvious benefits for applications other than parallel computing on distributed systems. For example, multimedia applications could benefit from coordinated scheduling of disk I/O, CPU cycles, and network bandwidth.

## References

- [Bir93] Kenneth Birman. Congestion control in Isis. Private communication, May 1993.
- [BJ87] K. Birman and T.A. Joseph. Exploiting virtual synchrony in distributed systems. In *Proc. 11th ACM Symp. on Operating System Principles, Austin, Texas*, pages 123–138, Nov. 1987.
- [BSS91] Kenneth Birman, André Schiper, and Pat Stephenson. Lightweight causal and atomic group multicast. *ACM Trans. on Computer Systems*, 9(3):272–314, Aug. 1991.
- [BT89] Dimitri P. Bertsekas and John N. Tsitsiklis. *Parallel and Distributed Computation: Numerical Methods*. Prentice Hall, Englewood Cliffs, New Jersey, 1989.
- [CBZ91] John B. Carter, John K. Bennett, and Willy Zwaenopoel. Implementation and performance of Munin. In *Proc. 13th ACM Symp. on Operating System Principles, Asilomar, California*, pages 152–164, Oct. 1991.
- [GK80] M. Gerla and L. Kleinrock. Flow control: a comparative survey. *IEEE Trans. Comm.*, COM-28:553–574, 1980.
- [HA90] P.W. Hutto and M. Ahamad. Slow memory: weakening consistency to enhance concurrency in distributed shared memories. In *Proc. 10th IEEE Intl. Conference on Distributed Computing Systems, Paris, France*, June 1990.
- [HM92] M. Herlihy and J.E.B. Moss. Transactional Memory: architectural support for lock-free data structures. Technical Report CRL 92/07, Digital Equipment Corp., Cambridge Research Lab, Dec. 1992.

- [HS92] Abdelsalam Heddaya and Himanshu S. Sinha. Coherence, non-coherence and Local Consistency in distributed shared memory for parallel computing. Technical Report BU-CS-92-004, Boston University, Computer Science Dept., May 1992.
- [HS93] Abdelsalam Heddaya and Himanshu S. Sinha. An overview of MERMERA: a system and formalism for non-coherent distributed parallel memory. In *Proc. 26th Hawaii International Conference on System Sciences, Maui, Hawaii*, pages 164–173, Jan. 5–8 1993.
- [HW91] Z. Haas and J. Winters. Congestion control by adaptive admission. In *Proc. IEEE INFOCOM*, pages 560–569, 1991.
- [Lam79] Leslie Lamport. How to make a multiprocessor computer that correctly executes multiprocess programs. *IEEE Trans. on Computers*, C-28(9):690–691, Sep. 1979.
- [LLJ+93] D. Lenoski, J. Laudon, T. Joe, D. Nakahira, L. Stevens, A. Gupta, and J. Hennessy. The DASH prototype: logic overhead and performance. *IEEE Trans. on Parallel and Distributed Systems*, 4(1):41–61, Jan. 1993.
- [LS88] R.J. Lipton and J.S. Sandberg. PRAM: a scalable shared memory. Technical Report CS-TR-180-88, Princeton Univ., Dept. of Computer Science, Sep. 1988.
- [Par93] Kihong Park. Warp control: a dynamically stable congestion protocol and its analysis. In *Proc. ACM SIGCOMM '93*, pages 137–147, September 1993.
- [Sin93] Himanshu Sinha. *MERMERA: Non-coherent Distributed Shared Memory for Parallel Computing*. PhD thesis, Boston University, Computer Science Department, 111 Cummington Street, Boston, MA 02215, May 1993.
- [Tan88] Andrew S. Tanenbaum. *Computer Networks*. Prentice-Hall, Englewood Cliffs, New Jersey, second edition, 1988.