

Sequential Consistency versus Linearizability

HAGIT ATTIYA
The Technion

and

JENNIFER L. WELCH
University of North Carolina at Chapel Hill

The power of two well-known consistency conditions for shared-memory multiprocessors, *sequential consistency* and *linearizability*, is compared. The cost measure studied is the worst-case response time in distributed implementations of virtual shared memory supporting one of the two conditions. Three types of shared-memory objects are considered: read/write objects, FIFO queues, and stacks. If clocks are only approximately synchronized (or do not exist), then for all three object types it is shown that linearizability is more expensive than sequential consistency: We present upper bounds for sequential consistency and larger lower bounds for linearizability. We show that, for all three data types, the worst-case response time is very sensitive to the assumptions that are made about the timing information available to the system. Under the strong assumption that processes have perfectly synchronized clocks, it is shown that sequential consistency and linearizability are equally costly: We present upper bounds for linearizability and matching lower bounds for sequential consistency. The upper bounds are shown by presenting algorithms that use atomic broadcast in a modular fashion. The lower-bound proofs for the approximate case use the technique of “shifting,” first introduced for studying the clock synchronization problem.

Categories and Subject Descriptors: B.3.2 [**Memory Structures**]: Design Styles—*cache memories*; *shared memory*; C.1.2 [**Processor Architectures**]: Multiple Data Stream Architectures—

This paper combines and unifies results that appear in preliminary form in Attiya and Welch 1991 and Attiya 1991. H. Attiya’s research was supported in part by the B. and G. Greenberg Research Fund (Ottawa), by Technion V.P.R. funds, and by the fund for the promotion of research at the Technion. J. L. Welch’s research was supported in part by NSF grant CCR-9010730, an IBM Faculty Development Award, and an NSF Presidential Young Investigator Award. Part of the work by Attiya was performed while visiting the DEC Cambridge Research Laboratory and the Laboratory for Computer Science, MIT, supported by ONR contract N00014-85-K-0168, by NSF grants CCR-8611442 and CCR-8915206, and by DARPA contracts N00014-89-J-1988 and N00014-87-K-0825.

Authors’ addresses: H. Attiya, Department of Computer Science, The Technion, Haifa 32000, Israel; J. L. Welch, Department of Computer Science, Texas A & M University, College Station, TX 77843-3112.

This paper was originally submitted to ACM TOPLAS, where it was reviewed and accepted for publication. It was subsequently transferred to TOCS, with the agreement of the authors, to take advantage of the shorter delay to publication for TOCS.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

© 1994 ACM 0164-0925/94/0500-0091\$03.50

parallel processors; D.3.3 [**Programming Techniques**]: Language Constructs—*abstract data types*; *concurrent programming structures*; D.4.2 [**Operating Systems**]: Storage Management—*distributed memories*; F.1.2 [**Computation by Abstract Devices**]: Modes of Computation—*parallelism and concurrency*

General Terms: Algorithms, Design

Additional Key Words and Phrases: Cache coherence, linearizability, sequential consistency, shared memory

1. INTRODUCTION

1.1 Overview

We present a quantitative comparison of the inherent costs of two well-known consistency conditions for concurrently accessed shared data: sequential consistency and linearizability. Our main conclusion is that under realistic timing assumptions it is strictly more expensive to provide linearizability than sequential consistency in distributed (message-passing) implementations, where the cost measure is the worst-case time for a data access. Because their definitions are very similar, linearizability and sequential consistency are often confused, but our work shows that there can be a considerable cost to such confusion. For example, Dubois and Scheurich [1990; Scheurich and Dubois 1987] define a sufficient condition for sequential consistency (cf. Dubois and Scheurich [1990, def. 6.3]); however, this condition implies linearizability.¹ Implementing linearizability imposes more cost than is necessary to support the target condition of sequential consistency. To our knowledge, this is the first time sequential consistency is shown to be more costly than linearizability.

We also study the worst-case access time for the two conditions under more stringent timing assumptions, namely, when processes have perfectly synchronized clocks. We show several lower bounds for sequential consistency in this model; these lower bounds carry over to more realistic models in which weaker assumptions hold about the clock behavior. We also present matching “counterexample” algorithms for linearizability in this model. They demonstrate that no improved lower bounds for either condition are possible without weakening the timing assumptions. The fact that sequential consistency and linearizability are equally costly in this model (for our measure) is somewhat surprising. It indicates the importance of explicitly and carefully specifying system timing assumptions, and the nontriviality of separating sequential consistency from linearizability.

1.2 Detailed Description

Managing concurrent accesses to shared data by several processes is a problem that arises in many contexts, ranging from cache coherence for

¹Technically, the definition of Dubois and Scheurich [1990] relies on the notion of “performing an operation,” which can only be interpreted in a specific architectural model. Under a natural interpretation (e.g., as in Gibbons et al. 1991), the definition implies linearizability.

multiprocessors to distributed file systems and transaction systems. A consistency condition must specify what guarantees are provided about the values returned by data accesses in the presence of interleaved and/or overlapping accesses, perhaps to distinct copies of a single logical data item. Two conflicting goals of a consistency condition are to be strong enough to be useful to the user and weak enough to be implemented efficiently. Sequential consistency and linearizability are two well-known consistency conditions.

Sequential consistency requires that all of the data operations *appear* to have executed atomically, in some sequential order that is consistent with the order seen at individual processes.² When this order must also preserve the global (external) ordering of nonoverlapping operations, this consistency guarantee is called *linearizability* [Herlihy and Wing 1990].³

Clearly, linearizability is stronger than sequential consistency. As discussed in Herlihy and Wing [1990], linearizability has two advantages over sequential consistency. First, it is somewhat more convenient to use, especially for formal verification, because it preserves real-time ordering of operations, and hence corresponds more naturally to the intuitive notion of atomic execution of operations. Second, linearizability is compositional (or local); that is, the combination of separate linearizable implementations of two objects yields a linearizable implementation.⁴ In contrast, sequential consistency is not compositional, implying that all objects must be implemented together. Consequently, development costs and the amount of synchronization needed increase, and it is harder to apply separate optimizations to different objects.

Several papers have proposed sequentially consistent implementations of read/write objects, which were claimed to achieve a high degree of concurrency (e.g., Adve and Hill [1990a], Afek et al. [1993], Bennett et al. [1990], Brantley et al. [1985], Dubois and Scheurich [1990], Min and Baer [1989], and Scheurich and Dubois [1987]). None of these papers proves that similar improvements cannot be achieved for linearizability, and none provides an analysis of the response time of the implementations suggested (or any other complexity measure).

In contrast, our work shows the existence of a gap between the fastest implementation of linearizability and what can be achieved for sequential consistency. To our knowledge, this is the first such quantitative result comparing the two conditions. Our results are shown in relatively abstract formal models. We believe that the correct abstraction has been done so that the results are applicable in a wide variety of contexts.

Our system model consists of a collection of application processes running concurrently and communicating via virtual shared memory. The shared

²This condition is similar to *serializability* from database theory [Bernstein et al. 1987; Papadimitriou 1986]; however, serializability applies to *transactions* that aggregate many operations.

³Also called *atomicity* [Herlihy 1988; Lamport 1986; Misra 1986] in the case of read/write objects.

⁴We use the term *implementation* in its usual meaning in the semantics literature, of satisfying a certain specification.

memory consists of a collection of *objects*. Unlike most previous research, we consider other types of shared objects in addition to the usual read/write objects. Since read/write objects do not provide an expressive and a convenient abstraction for concurrent programming (cf. Herlihy [1988]), many multiprocessors now support more powerful concurrent objects, for example, FIFO queues, stacks, Test & Set, and Fetch & Add [Brantley et al. 1985; Gottlieb et al. 1983]. We study FIFO queues and stacks; our results easily extend to Test & Set and Fetch & Add.

The application processes are running in a distributed system consisting of a collection of nodes and a communication network. The network need not be fully connected physically, but it must be possible for every node to communicate with every other node. The shared-memory abstraction is implemented by a *memory-consistency system (MCS)*, which uses local memory at the various nodes and some protocol executed by the MCS processes (one at each node). (Nodes that are dedicated storage can be modeled by nullifying the application process.) Figure 1 illustrates a node, on which an application process and an MCS process are running. The application process sends calls to access shared data to the MCS process; the MCS process returns the responses to the application process, possibly based on messages exchanged with MCS processes on other nodes. The responses must be consistent with the particular consistency condition that is being implemented. Thus, the consistency condition is defined at the interface between the application process and the MCS.

On each node there is a real-time clock readable by the MCS process at that node, which runs at the same rate as real time. Every message incurs a delay in the interval $[d - u, d]$, for some known constants u and d , $0 \leq u \leq d$ (u stands for *uncertainty*). If $u = 0$, then the message delays are constant. Thus, d is the worst-case delay in the network over all pairs of processes, and u is the worst-case uncertainty. Our lower bounds are given in terms of d and u . Our upper bounds are given in terms of d (u happens not to enter in); however, all of our algorithms also work in asynchronous systems where d and u are unknown (or even unbounded), since these constants appear nowhere in the code.

We have chosen to focus on the distinction between performing a data operation locally at a process, based on its local state, and performing an operation that requires communication between processes before it can return to the user. We model this by assuming 0 time for local processing and $d > 0$ time for the worst-case communication cost in the system. This is a reasonable approximation, as in many systems the time required for communication far outweighs the time for local computation.

We start with the case when the process clocks are only approximately synchronized and there is uncertainty in the message delay.⁵ Under this assumption, for all three object types, there are gaps between the upper bounds for sequentially consistent implementations and the lower bounds for

⁵If there is no uncertainty in the message delay, then clocks can be perfectly synchronized.

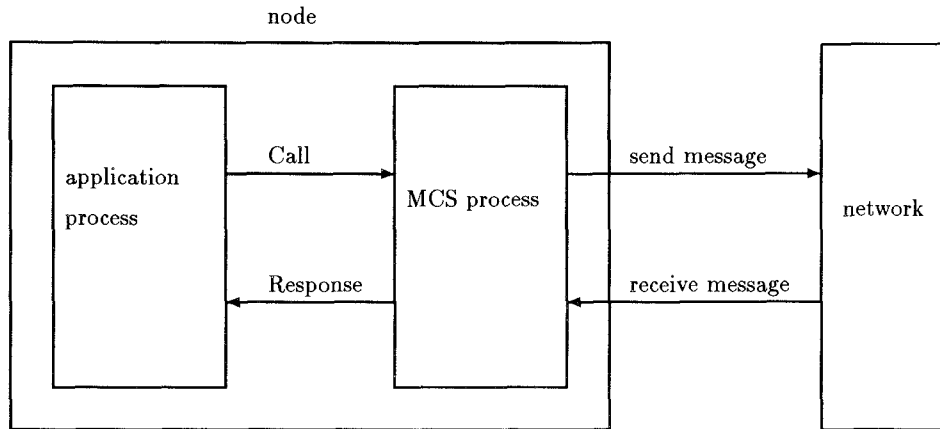


Fig. 1. System architecture.

linearizable implementations. We show that there are operations that can be done instantaneously (i.e., locally) in sequentially consistent implementations, but that require $\Omega(u)$ time in linearizable implementations (note that u can be as large as d). In particular,

- for read/write objects,
 - for linearizability, the worst-case time for a read is at least $u/4$;
 - for linearizability, the worst-case time for a write is at least $u/2$; and
 - for sequential consistency, there is an algorithm that guarantees time 0 for a read and time $2d$ for a write, and another that guarantees the reverse;
- for FIFO queues,
 - for linearizability, the worst-case time for an enqueue is at least $u/2$; and
 - for sequential consistency, there is an algorithm that guarantees time 0 for an enqueue and time $2d$ for a dequeue;
- the situation for stacks is analogous to that for FIFO queues, with “pop” playing the role of “dequeue” and “push” the role of “enqueue.”

Thus, under these timing assumptions, linearizability is more expensive to implement than sequential consistency, when there are significantly more operations of one type.

We then consider the stronger model when processes’ clocks are perfectly synchronized. In this case, we show several strong lower bounds for sequential consistency. We also give matching upper bounds for linearizability. In particular,

- for read/write objects,
 - for sequential consistency, the sum of the worst-case response times for a read operation and a write operation is at least d ; and

- for linearizability, there is an algorithm in which a read operation is performed instantaneously (locally), while a write operation returns within time d ; also, there is an algorithm in which the roles are reversed.

The above lower bound formalizes and strengthens a result of Lipton and Sandberg [1988]. The two matching algorithms show that the lower-bound trade-off is tight: It is possible to have the response time of *only one* of the operations depend on the network's latency.

- for FIFO queues,
 - for sequential consistency, the worst-case response time for a dequeue operation is at least d ; and
 - for linearizability, there is an algorithm in which an enqueue operation returns instantaneously, while a dequeue operation returns within time d ;
- for stacks, as in the case of imperfect clocks, the results are analogous to those for FIFO queues, with “pop” playing the role of “dequeue” and “push” the role of “enqueue.”

Thus, we need to assume that clocks are imperfect in order to separate sequential consistency from linearizability, indicating that this separation is not as obvious as it may seem and depends on delicate timing assumptions.

Section 2 presents our definitions and reviews a technique called “shifting” that is used in our lower bounds. Section 3 covers the imperfect-clock case. There is a subsection for each object type; each subsection consists of the lower bound(s) for linearizability, followed by the upper bound(s) for sequential consistency. Section 4 considers the perfect-clock case. Again, there is one subsection for each of the three object types; each subsection consists of the lower bound(s) for sequential consistency followed by the upper bound(s) for linearizability. We conclude in Section 5 with a discussion of our results, describe work that followed the original version of this paper, and suggest avenues for further research.

2. PRELIMINARIES

2.1 Objects

Every shared object is assumed to have a *serial specification* (cf., Herlihy and Wing [1990]) defining a set of *operations*, which are ordered pairs of call and response events, and a set of operation sequences, which are the allowable sequences of operations on that object. A sequence τ of operations for a collection of objects is *legal* if, for each object O , the restriction of τ to operations of O , denoted $\tau|O$, is in the serial specification of O .

In the case of a read/write object X , the ordered pair of events $[\text{Read}_p(X), \text{Ret}_p(X, v)]$ forms a Read operation for any process p and value v , and $[\text{Write}_p(X, v), \text{Ack}_p(X)]$ forms a Write operation. The set of operation sequences consists of all sequences in which every read operation returns the

value of the latest preceding write operation (the usual read/write semantics).⁶

In the case of a FIFO queue Q , the ordered pair of events $[\text{Deq}_p(Q), \text{Ret}_p(Q, v)]$ forms a Deq operation for any process p and value v , and $[\text{Enq}_p(Q, v), \text{Ack}_p(Q)]$ forms an Enq operation. The set of operation sequences consists of all sequences that obey the usual FIFO queue semantics. That is, with a sequence of operations we associate a sequence of *queue states*, starting with an initial empty state and continuing with a state for each operation representing the state of the queue *after* the operation. Each enqueue operation adds an item to the end of the queue, and each dequeue operation removes an item from the head of the queue, or returns \perp if the queue is empty.

The specification of a stack S is similar to the specification of a queue: $[\text{Pop}_p(S), \text{Ret}_p(S, v)]$ forms a Pop operation for any process p and value v , and $[\text{Push}_p(S, v), \text{Ack}_p(S)]$ forms a Push operation. The set of operation sequences consists of all sequences that obey the usual last-in-first-out stack semantics.

2.2 System Model

We assume a system consisting of a collection of nodes connected via a communication network. On each node there is an application process, an MCS process, and a real-time clock readable by the MCS process at that node. Formally, a *clock* is a monotonically increasing function from \mathcal{R} (real time) to \mathcal{C} (clock time).⁷ The clock cannot be modified by the process. Processes do not have access to the real time; each process obtains its only information about time from its clock.

Below we list and informally explain the *events* that can occur at the MCS process on node p . (The name p is also used for the MCS process on node p .)

Call events. The application process on node p wants to access a shared object.

Response events. The MCS process on node p is providing a response from a shared object to the application process on node p .

Message receive events. $\text{receive}(p, m, q)$ for all messages m and nodes q : The MCS process on node p receives message m from the MCS process on node q .

Message send events. $\text{send}(p, m, q)$ for all messages m and MCS processes q : The MCS process on node p sends message m to the MCS process on node q .

Timer-set events. $\text{timerset}(p, T)$ for all clock times T : p sets a timer to go off when its clock reads T .

⁶The specifications used in this paper are operational. It is possible to give algebraic (axiomatic) specifications (cf. [Herlihy and Wing 1990]); operational specifications are used here for simplicity.

⁷ \mathcal{R} denotes the real numbers.

Timer events. $\text{timer}(p, T)$ for all clock times T : A timer that was set for time T on p 's clock goes off.

The call, message-receive, and timer events are *interrupt* events.

An *MCS process* (or simply *process*) is an automaton with a (possibly infinite) set of states, including an initial state, and a transition function. Each interrupt event causes an application of the transition function. The transition function is a function from states, clock times, and interrupt events to states, sets of response events, sets of message-send events, and sets of timer-set events (for subsequent clock times). That is, the transition function takes as input the current state, clock time, and interrupt event (which is the receipt of a call from the application process, the receipt of a message from another node, or a timer going off), and produces a new state, a set of response events for the application process, a set of messages to be sent and a set of timers to be set for the future.

A *step* of p is a tuple (s, T, i, s', R, M, S) where s and s' are states, T is a clock time, i is an interrupt event, R is a set of response events, M is a set of message-send events, S is a set of timer-set events, and $s', R, M,$ and S are the result of p 's transition function acting on $s, T,$ and i .

A *history* of a process p with clock C is a mapping from \mathcal{R} (real time) to finite sequences of steps such that

- (1) for each real time t , there is only a finite number of times $t' < t$ such that the corresponding sequence of steps is nonempty (thus, the concatenation of all the sequences in real-time order is a sequence);
- (2) the old state in the first step is p 's initial state;
- (3) the old state of each subsequent step is the new state of the previous step;
- (4) for each real time t , the clock time component of every step in the corresponding sequence is equal to $C(t)$; and
- (5) for each real time t , in the corresponding sequence all nontimer events are ordered before any timer event, and there is at most one timer event.

A *memory-consistency system (MCS)* is a set of processes P together with a set of clocks \mathcal{C} , one for each p in P . An *execution* of an MCS is a set of histories, one for each process p in P with clock C_p in \mathcal{C} , satisfying the following two conditions: (1) There is a one-to-one correspondence between the messages sent by p to q and the messages received by q from p , for any processes p and q . We use the message correspondence to define the *delay* of any message in an execution to be the real time of receipt minus the real time of sending. (2) A timer is received by p at clock time T iff p has previously set a timer for T . (The network is not explicitly modeled, although the constraints on executions given below imply that the network reliably delivers all messages sent.)

Execution σ is *admissible* if the following conditions hold:

- (1) For every p and q , every message in σ from p to q has its delay in the range $[d - u, d]$, for fixed nonnegative integers d and $u, u \leq d$. (This is a restriction on the network.)

- (2) For every p , at most one call at p is pending at a time. (This is a restriction on the application process.)

Note that the last condition allows each application process to have at most one call outstanding at any time. This outlaws pipelining or prefetching.

2.3 Correctness Conditions

Given an execution σ , let $ops(\sigma)$ be the sequence of call and response events appearing in σ in real-time order, breaking ties for each real time t as follows: First, order all response events for time t whose matching call events occur before time t , using process IDs to break any remaining ties. Then order all operations whose call and response both occur at time t . Preserve the relative ordering of operations for each process, and break any remaining ties with process IDs. Finally, order all call events for time t whose matching response events occur after time t , using process IDs to break any remaining ties.

Our formal definitions of sequential consistency and linearizability follow. These definitions imply that every call gets an eventual response and that calls and responses alternate at each process. Given a sequence s of operations and a process p , we denote by $s|p$ the restriction of s to operations of p .

Definition 2.1 (Sequential consistency). An execution σ is *sequentially consistent* if there exists a legal sequence τ of operations such that τ is a permutation of $ops(\sigma)$ and, for each process p , $ops(\sigma)|p$ is equal to $\tau|p$.

Definition 2.2 (Linearizability). An execution σ is *linearizable* if there exists a legal sequence τ of operations such that τ is a permutation of $ops(\sigma)$, for each process p , $ops(\sigma)|p$ is equal to $\tau|p$, and furthermore, whenever the response for operation op_1 precedes the call for operation op_2 in $ops(\sigma)$, then op_1 precedes op_2 in τ .

An MCS is a *sequentially consistent* implementation of a set of objects if any admissible execution of the MCS is sequentially consistent; similarly, an MCS is a *linearizable* implementation of a set of objects if any admissible execution of the MCS is linearizable.

We measure the efficiency of an implementation by the worst-case response time for any operation on any object in the set. Given a particular MCS, an object O implemented by it, and an operation P on O , we denote by $|P(O)|$ the maximum time taken by a P operation on O in any admissible execution. We denote by $|P|$ the maximum of $|P(O)|$ over all objects O implemented by the MCS.

2.4 Shifting

A basic technique we use in our lower-bound proofs (in Section 3.1.1 and 3.2.1) is *shifting*, originally introduced in Lundelius and Lynch [1984] to prove lower bounds on the precision achieved by clock synchronization

algorithms. Shifting is used to change the timing and the ordering of events in the system, while preserving the local views of the processes.

Informally, given an execution with a certain set of clocks, if process p 's history is changed so that the real times at which the events occur are shifted by some amount s and if p 's clock is shifted by the same amount, then the result is another execution in which every process still “sees” the same events happening at the same clock times. The intuition is that the changes in the real times at which events happen at p cannot be detected by p because its clock has changed by a corresponding amount.

More precisely, the *view* of process p in history π of p with clock C is the concatenation of the sequences of steps in π , in real-time order. The real times of occurrence are not represented in the view. Two histories, one of process p with clock C and the other of process p with clock C' , are *equivalent* if the view of p is the same in both histories. Two executions, execution σ of system (P, \mathcal{C}) and execution σ' of (P, \mathcal{C}') , are *equivalent* if, for each process p , the component histories for p in σ and σ' are equivalent. Thus, the executions are indistinguishable to the processes. Only an outside observer who has access to the real time can tell them apart.

Given history π of process p with clock C , and real number s , a new history $\pi' = \text{shift}(\pi, s)$ is defined by $\pi'(t) = \pi(t + s)$ for all t . That is, all tuples are shifted earlier in π' by s if s is positive, and later by $|s|$ if s is negative. Given a clock C and real number s , a new clock $C' = \text{shift}(C, s)$ is defined by $C'(t) = C(t) + s$ for all t . That is, the clock is shifted forward by s if s is positive, and backward by $|s|$ if s is negative.

The following lemma observes that shifting a history of process p and p 's clock by the same amount produces another history:

LEMMA 2.1. *Let π be a history of process p with clock C , and let s be a real number. Then $\text{shift}(\pi, s)$ is a history of p with clock $\text{shift}(C, s)$.*

Given execution σ of system (P, \mathcal{C}) and real number s , a new execution $\sigma' = \text{shift}(\sigma, p, s)$ is defined by replacing π , p 's history in σ , by $\text{shift}(\pi, s)$, and by retaining the same correspondence between sends and receives of messages. (Technically, the correspondence is redefined so that a pairing in σ that involves the event for p at time t , in σ' , involves the event for p at time $t - s$.) All tuples for process p are shifted by s , but no others are altered. Given a set of clocks $\mathcal{C} = \{C_q : q \in P\}$, and real number s , a new set of clocks $\mathcal{C}' = \text{shift}(\mathcal{C}, p, s)$ is defined by replacing clock C_p by clock $\text{shift}(C_p, s)$. Process p 's clock is shifted forward by s , but no other clocks are altered.

The following lemma observes that shifting one process's history and clock by the same amount in an execution results in another execution that is equivalent to the original:

LEMMA 2.2 [Lundelius and Lynch 1984]. *Let σ be an execution of system (P, \mathcal{C}) , p a process, and s a real number. Let $\mathcal{C}' = \text{shift}(\mathcal{C}, p, s)$ and $\sigma' = \text{shift}(\sigma, p, s)$. Then σ' is an execution of (P, \mathcal{C}') , and σ' is equivalent to σ .*

The following lemma quantifies how message delays change when an execution is shifted. Notice that the result of shifting an admissible execution is not necessarily admissible:

LEMMA 2.3 [Lundelius and Lynch 1984]. *Let σ be an execution of system (P, \mathcal{E}) , p a process, and s a real number. Let $\mathcal{E}' = \text{shift}(\mathcal{E}, p, s)$ and $\sigma' = \text{shift}(\sigma, p, s)$. Make the obvious correspondence between messages in σ and in σ' . Suppose x is the delay of message m from process q to process r in σ . Then the delay of m in σ' is x if $q \neq p$ and $r \neq p$, $x - s$ if $r = p$, and $x + s$ if $q = p$.*

3. IMPERFECT CLOCKS

We start by assuming a system in which clocks run at the same rate as real time but are not initially synchronized, and in which message delays are in the range $[d - u, d]$ for some $u > 0$. We show that in this model there is a gap between the upper bounds for sequential consistency and the lower bounds for linearizability, for all three object types. The lower bounds state that the worst-case time for a read is at least $u/4$ and the worst-case time for a write, an enqueue, or a push is at least $u/2$. Recall that u is the uncertainty in the message delay and can be as large as d . In contrast, we describe sequentially consistent algorithms that implement these operations instantaneously.

Intuitively, the algorithms are similar to a snoopy write-through cache protocol with a write-broadcast policy for bus-based systems. The main idea is that, in order to guarantee sequential consistency, it suffices for processes to update their local copies in the same order. (This is provided immediately in bus-based systems.) A simple way to achieve this property is for a centralized controller to collect update messages and to broadcast them. Using atomic broadcast it is possible to translate this idea into algorithms that are fully distributed and that do not rely on a centralized controller. The algorithms do not rely on timing information and also work in an asynchronous system.

Atomic broadcast [Birman and Joseph 1987] is a communication primitive that guarantees that every message sent using the primitive is received at every process, all messages are delivered in the same order at all processes, and two messages sent by the same process are delivered in the same order they were sent. Our implementations are described in a modular way so that they will work with any atomic broadcast algorithm (e.g., Birman and Joseph [1987], Chang and Maxemchuk [1984], and Garcia-Molina and Spauster [1989]). The interface to the primitive consists of two operations, ABC-send(m) to broadcast a message m (possibly consisting of several fields) and ABC-receive(m) to receive a message m . In analyzing our implementations, we assume there is a known bound, h , on the time that the atomic broadcast primitive takes to deliver a message to all processes. Each of our implementations has one fast operation, which takes time 0, and one slow operation, which takes time h . In the Appendix we describe and prove correct a fast atomic broadcast algorithm with $h = 2d$. By using this algorithm in our implementations, we obtain implementations in which slow operations take time $2d = O(d)$.

3.1 Read/Write Objects

We show in Section 3.1.1 that, in any linearizable implementation of a read/write object, the worst-case response time of *both* read and write operations must depend on u . We then present in Section 3.1.2 two sequentially consistent algorithms for read/write objects, one in which reads are performed instantaneously while the worst-case response time for a write is $O(d)$, and another in which the roles are reversed.

3.1.1 Lower Bounds for Linearizability. We now show that, under reasonable assumptions about the pattern of sharing, in any linearizable implementation of an object, the worst-case time for a read is $u/4$ and the worst-case time for a write is $u/2$. The proofs of these lower bounds use the technique of shifting, described in Section 2.4.

THEOREM 3.1. *Assume X is a read/write object with at least two readers and a distinct writer. Then any linearizable implementation of X must have $|Read(X)| \geq u/4$.*

PROOF. Let p and q be two processes that read X , and let r be a process that writes X . Assume in contradiction that there is an implementation with $|Read(X)| < u/4$. Without loss of generality, assume that the initial value of X is 0. The idea of the proof is to consider an execution in which p reads 0 from X , then q and p alternate reading X while r writes 1 to X , and then q reads 1 from X . Thus, there exists a read R_1 , say, by p , that returns 0 and is immediately followed by a read R_2 by q that returns 1. If q is shifted earlier by $u/2$, then R_2 precedes R_1 in the resulting execution. Since R_2 returns the new value 1 and R_1 returns the old value 0, this contradicts linearizability.

Let $k = \lceil |Write(X)|/u \rceil$. By the specification of X , there is an admissible execution α , in which all message delays are $d - (u/2)$, consisting of the following operations (see Figure 2a⁸):

- At time $u/4$, r does a $Write_r(X, 1)$.
- Between times $u/4$ and $(4k + 1) \cdot (u/4)$, r does an $Ack_r(X)$. (By definition of k , $(4k + 1) \cdot (u/4) \geq (u/4) + |Write(X)|$, and thus, r 's write operation is guaranteed to finish in this interval.)
- At time $2i \cdot (u/4)$, p does a $Read_p(X)$, $0 \leq i \leq 2k$.
- Between times $2i \cdot (u/4)$ and $(2i + 1) \cdot (u/4)$, p does a $Ret_p(X, v_{2i})$, $0 \leq i \leq 2k$.
- At time $(2i + 1) \cdot (u/4)$, q does a $Read_q(X)$, $0 \leq i \leq 2k$.
- Between times $(2i + 1) \cdot (u/4)$ and $(2i + 2) \cdot (u/4)$, q does a $Ret_q(X, v_{2i+1})$, $0 \leq i \leq 2k$.

Thus, in $ops(\alpha)$, p 's read of v_0 precedes r 's write, q 's read of v_{4k+1} follows r 's write, no two read operations overlap, and the order of the values read

⁸In the figures, time runs from left to right, and each line represents events at one process. Important time points are marked at the bottom.

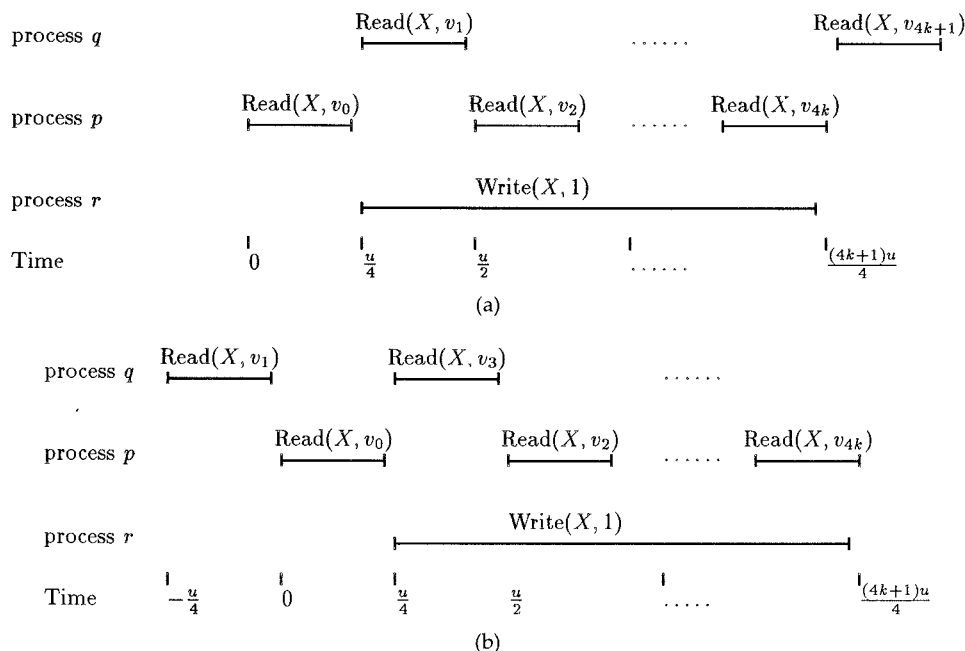


Fig. 2. Executions used in the proof of Theorem 3.1 (a) The execution α . (b) The execution β .

from X is $v_0, v_1, v_2, \dots, v_{4k+1}$. By linearizability, $v_0 = 0$ and $v_{4k+1} = 1$. Thus, there exists $j, 0 \leq j \leq 4k$, such that $v_j = 0$ and $v_{j+1} = 1$. Without loss of generality, assume that j is even, so that v_j is the result of a read by p .

Define $\beta = \text{shift}(\alpha, q, (u/2))$; that is, we shift q earlier by $u/2$ (see Figure 2b). The result is admissible, since by Lemma 2.3 the message delays to q become $d - u$, the message delays from q become d , and the remaining message delays are unchanged.

As a result of the shifting, we have reordered read operations with respect to each other at p and q . Specifically, in $\text{ops}(\beta)$, the order of the values read from X is $v_1, v_0, v_3, v_2, \dots, v_{j+1}, v_j, \dots$. Thus, in β we now have $v_{j+1} = 1$ being read before $v_j = 0$, which violates linearizability. \square

THEOREM 3.2. *If X is a read/write object with at least two writers and a distinct reader, then any linearizable implementation of X must have $|\text{Write}(X)| \geq (u/2)$.*

The proof uses techniques similar to the proof of Theorem 3.1. It constructs an execution in which, if write operations are too short, linearizability can be violated by appropriately shifting histories.

PROOF. Let p and q be two processes that write X , and let r be a process that reads X . Assume in contradiction that there is an implementation with

$|\text{Write}(X)| < u/2$. Without loss of generality, assume that the initial value of X is 0. By the specification of X , there is an admissible execution α such that

- $\text{ops}(\alpha)$ is $\text{Write}_p(X, 1) \text{Ack}_p(X) \text{Write}_q(X, 2) \text{Ack}_q(X) \text{Read}_r(X) \text{Ret}_r(X, 2)$;
- $\text{Write}_p(X, 1)$ occurs at time 0, $\text{Write}_q(X, 2)$ occurs at time $u/2$, and $\text{Read}_r(X)$ occurs at time u ; and
- the message delays in α are d from p to q , $d - u$ from q to p , and $d - (u/2)$ for all other ordered pairs of processes.

Let $\beta = \text{shift}(\text{shift}(\alpha, p, -(u/2)), q, (u/2))$; that is, we shift p later by $u/2$ and q earlier by $u/2$. The result is still an admissible execution, since by Lemma 2.3 the delay of a message from p or to q becomes $d - u$, the delay of a message from q or to p becomes d , and all other delays are unchanged.

But $\text{ops}(\beta)$ is

$$\text{Write}_q(X, 2) \text{Ack}_q(X) \text{Write}_p(X, 1) \text{Ack}_p(X) \text{Read}_r(X) \text{Ret}_r(X, 2),$$

which violates linearizability, because r 's read should return 1, not 2. \square

The assumptions about the number of readers and writers made in Theorems 3.1 and 3.2 are crucial to the results, since it can be shown that the algorithms from Theorems 4.2 and 4.3 are correct if there is only one reader and one writer.

3.1.2 Upper Bounds for Sequential Consistency

Fast reads. We start with the algorithm for fast reads (time 0) and slow writes (time at most h).

In the algorithm, each process keeps a local copy of every object. A read returns the value of the local copy immediately. When a write comes into p , p sends an atomic broadcast containing the name of the object to be updated and the value to be written; but it does not yet generate an Ack for the write operation. When an update message is delivered to a process q , q writes the new value to its local copy of the object. If the update message was originated by q , then q generates an Ack, and the (unique pending) write operation returns.

More precisely, the state of each process consists of a copy of every object, initially equal to its initial value. The transition function of process p appears in Figure 3.

To prove the correctness of the algorithm, we first show the following lemma:

LEMMA 3.3. *For every admissible execution and every process p , p 's local copies take on all of the values contained in write operations, all updates occur in the same order at each process, and this order preserves the order of write operations on a per-process basis.*

PROOF. By the code, an ABC-send is done exactly once for each write operation. By the guarantees of the atomic broadcast, each process receives exactly one message for each write operation, these messages are received in

```

Readp(X):
  generate Retp(X, v), where v is the value of p's copy of X

Writep(X, v):
  ABC-send(X, v)

ABC-receive(X, v) from q:
  set local copy of X to v
  if q = p then generate Ackp(X) endif

```

Fig. 3. Sequentially consistent fast-read algorithm.

the same order at each process, and this order respects the order of sending on a per-process basis. \square

Call the total order of Lemma 3.3 the “Abcast order.”

LEMMA 3.4. *For every admissible execution, every process p , and all objects X and Y , if read R of object Y follows write W to object X in $ops(\sigma)|p$, then R 's read of p 's local copy of Y follows W 's write of p 's local copy of X .*

PROOF. The lemma is true because W does not end until its update is performed at its initiator. \square

THEOREM 3.5. *There exists a sequentially consistent implementation of read/write objects with $|Read| = 0$ and $|Write| = h$.*

PROOF. Consider the algorithm just presented. Clearly, the time for any read is 0. The time for any write is the time for the initiator's ABC-send to be received by the initiator, which is at most h . The remainder of the proof is devoted to showing sequential consistency. Fix admissible execution σ .

Define the sequence of operations τ as follows: Order the writes in σ in Abcast order. Now we explain where to insert the reads. We proceed in order from the beginning of σ . $[Read_p(X), Ret_p(X, v)]$ goes immediately after the latest of (1) the previous operation for p (either read or write, on any object) and (2) the write that spawned the latest update of p 's local copy of X preceding the generation of the $Ret_p(X, v)$. (Break ties using process IDs; e.g., if every process reads some object before any process writes any object, then τ begins with p_1 's read, followed by p_2 's read, etc.)

We must show $ops(\sigma)|p = \tau|p$ for all processes p . Fix some process p . The relative ordering of two reads in $ops(\sigma)|p$ is the same in $\tau|p$ by definition of τ . The relative ordering of two writes in $ops(\sigma)|p$ is the same in $\tau|p$ by Lemma 3.3. Suppose in $ops(\sigma)|p$ that read R follows write W . By definition of τ , R comes after W in τ .

Suppose in $ops(\sigma)|p$ that read R precedes write W . Suppose in contradiction that R comes after W in τ . Then in σ there is some read $R' = [Read_p(X), Ret_p(X, v)]$ and some write $W' = [Write_q(x, v), Ack_q(X)]$ such that (1) R' equals R or occurs before R in σ , (2) W' equals W or follows W in the Abcast order, and (3) W' spawns the latest update to p 's copy of X that

precedes R 's read. But, in σ , R finishes before W starts. Since updates are performed in σ in Abcast order (Lemma 3.3), R cannot see W 's update, a contradiction.

We must show that τ is legal. Consider read $R = [\text{Read}_p(X), \text{Ret}_p(X, v)]$ in τ . Let W be the write in σ that spawns the latest update to p 's copy of X preceding R 's read of p 's copy of X . Clearly, $W = [\text{Write}_q(X, v), \text{Ack}_q(X)]$ for some q . (If there is no such W , then consider an imaginary write at the beginning of σ .) By the definition of τ , R follows W in τ . We must show that no other write to X falls in between W and R in τ . Suppose in contradiction that $W' = [\text{Write}_r(X, w), \text{Ack}_r(X)]$ does. Then, by Lemma 3.3, the update for W' follows the update for W at every process in σ .

Case 1: $r = p$. Since τ preserves the order of operations at p , W' precedes R in σ . Since the update for W' follows the update for W in σ , R sees W' 's update, not W 's, contradicting the choice of W .

Case 2: $r \neq p$. By definition of τ , there is some operation in $\text{ops}(\sigma)|p$ that, in τ , precedes R and follows W' (otherwise, R would not follow W'). Let O be the first such operation.

Suppose O is a write to some object Y . By Lemma 3.4, O 's update to p 's copy of Y precedes R 's read of p 's copy of X . Since updates are done in Abcast order, the update for W' occurs at p before the update for O , and thus, before R 's read, contradicting the choice of W .

Suppose O is a read. By the definition of τ , O is a read of X , and W' 's update to p 's copy of X is the latest one preceding O 's read (otherwise, O would not follow W'). Since updates are done in Abcast order, the value from W' supersedes the value from W , contradicting the choice of W . \square

Theorem 3.1 implies that this algorithm does not guarantee linearizability. We can also explicitly construct an admissible execution that violates linearizability as follows: The initial value of X is 0. Process p writes 1 to X . The ABC-send for the write occurs at time t . It arrives at process r at time t and at process q at time $t + h$. Meanwhile, r performs a read at time t and gets the new value 1, while q performs a read at time $t + h/2$ and gets the old value 0. No permutation of these operations can both conform to the read/write specification and preserve the relative real-time orderings of all nonoverlapping operations.

Fast writes. We now discuss the algorithm that ensures sequential consistency with fast writes (time 0) and slow reads (time at most h). When a $\text{Read}(X)$ comes into p , if p has no pending updates (to any object, not just X) that it initiated, then it returns the current value of its copy of X . Otherwise, it waits for all pending writes to complete and then returns. This is done by maintaining a count of the pending writes and waiting for it to be zero. When a $\text{Write}(X)$ comes into p , it is handled very similarly to the other algorithm; however, it is Acked immediately. Effectively, the algorithm

pipelines write updates generated at the same process. Specifically, the state of each process consists of the following variables:

- num: integer, initially 0 (number of pending updates initiated by this process); and
- copy of every object, initially equal to its initial value.

The transition function of process p appears in Figure 4.

THEOREM 3.6. *There exists a sequentially consistent implementation of read/write objects with $|Read| = h$ and $|Write| = 0$.*

PROOF. Consider the algorithm just presented. Clearly, every write takes 0 time. The worst-case time for a read occurs if the return must wait for the initiator to receive its own ABC-send for a pending write. This takes at most h time. The structure of the proof of sequential consistency is identical to that in the proof of Theorem 3.5. We just need a new proof for Lemma 3.4. Lemma 3.4 is still true for this algorithm because when a Read occurs at p , if any update initiated by p is still waiting, then the Return is delayed until the latest such update is performed. \square

Theorem 3.2 implies that this algorithm does not guarantee linearizability. An explicit scenario is easy to construct as well.

3.2 FIFO Queues

We show in Section 3.2.1 that, in any linearizable implementation of a FIFO queue, the worst-case response time of an enqueue operation must depend on u . We then present in Section 3.2.2 a sequentially consistent implementation in which enqueue operations return instantaneously, while the worst-case response time for a dequeue operation is h .

3.2.1 Lower Bound for Linearizability. We show that in any linearizable implementation of a FIFO queue the worst-case time for an enqueue is $u/2$ (assuming that at least two processes can enqueue to the same FIFO queue). The proof uses the technique of shifting, described in Section 2.4.

THEOREM 3.7. *If Q is a FIFO queue with at least two enqueueers and a distinct dequeuer, then any linearizable implementation of Q must have $|Enq(Q)| \geq (u/2)$.*

PROOF. Let p and q be two processes that can enqueue to Q , and let r be a process that dequeues from Q . Assume in contradiction that there is an implementation with $|Enq(Q)| < (u/2)$. Initially, Q is empty. By the specification of Q , there is an admissible execution α such that

- $ops(\alpha)$ is $Enq_p(Q, 1) Ack_p(Q) Enq_q(Q, 2) Ack_q(Q) Deq_r(Q) Ret_r(Q, 1)$;
- $Enq_p(Q, 1)$ occurs at time 0, $Enq_q(Q, 2)$ occurs at time $u/2$, and $Deq_r(Q)$ occurs at time u ; and
- the message delays in α are d from p to q , $d - u$ from q to p , and $d - (u/2)$ for all other ordered pairs of processes.

```

Readp(X):
  if num = 0 then
    generate Retp(X, v), where v is the value of p's copy of X
  endif

Writep(X, v)
  num := num + 1
  ABC-send(X, v)
  generate Ackp(X)

ABC-receive(X, v, t) from q:
  set local copy of X to v
  if p = q then
    num := num - 1
    if num = 0 then
      generate Retp(X, v), where v is the value of p's copy of X
    endif
  endif
endif

```

Fig. 4. Sequentially consistent fast-write algorithm.

Let $\beta = \text{shift}(\text{shift}(\alpha, p, -(u/2)), q, u/2)$; that is, we shift p later by $u/2$ and q earlier by $u/2$. The result is still an admissible execution, since by Lemma 2.3 the delay of a message from p or to q becomes $d - u$, the delay of a message from q or to p becomes d , and all other delays are unchanged. But $\text{ops}(\beta)$ is

$$\text{Enq}_q(Q, 2) \text{Ack}_q(Q) \text{Enq}_p(Q, 1) \text{Ack}_p(Q) \text{Deq}_r(Q) \text{Ret}_r(Q, 1),$$

which violates linearizability, because r 's dequeue should return 2, not 1 (by the FIFO property). \square

The assumption about the number of enqueueers made in Theorem 3.7 is crucial to the results, since it can be shown that the algorithm of Theorem 4.6 is correct if there is only one enqueueer.

3.2.2 Upper Bound for Sequential Consistency. Informally, the algorithm works as follows: Each process keeps a local copy of every object. When a request to enqueue v to Q comes into p , p broadcasts an update message with the object name, the operation name, and the value to be enqueued to all processes. The operation returns immediately. When a request to dequeue from Q comes into p , p broadcasts an update message with the object name and the operation name. It does not generate a response.

When an update message (either “deq” or “enq”) is delivered to a process, it handles it by performing the appropriate change (enqueue or dequeue) to the local copy of the object. If the update is a dequeue by the same process, the dequeue operation that is currently waiting returns the value that was

dequeued from the local copy. (Note that, by well formedness, there is only one pending dequeue operation for a given process.)

In more detail, the state of each process consists of the following variables:

- copy of every object, initially equal to its initial value; and
- val: value (of a queue element).

The transition function of process p appears in Figure 5.

To prove correctness of the algorithm, we show the following lemma:

LEMMA 3.8. *In every admissible execution, all updates are done exactly once at each local copy, updates are done in the same order at each process, and this order preserves the per-process order.*

THEOREM 3.9. *There exists a sequentially consistent implementation of FIFO queues with $|Enq| = 0$ and $|Deq| = h$.*

PROOF. Consider the algorithm just presented. Clearly, the time for an enqueue is 0, and the time for a dequeue is at most h . The remainder of the proof is devoted to showing sequential consistency. Fix some admissible execution σ .

Define the sequence of operations τ as follows: Order the operations in σ by Abcast order. From Lemma 3.8 it follows that operations by p are ordered in τ as they were ordered in σ ; thus, $ops(\sigma)|_p = \tau|_p$, for all processes p . It remains to show that τ is legal, that is, that for every FIFO queue Q , $\tau|_Q$ is in the serial specification of Q . Pick any Q , and consider $\tau|_Q = op_1op_2\dots$. Suppose op_i is $[Deq_p(Q), Ret_p(Q, v)]$. Since the local updates at p occur in Abcast order (Lemma 3.8), updates at p to the local copy of Q occur in the same order as in τ , and the claim follows. \square

Theorem 3.7 implies that this algorithm does not guarantee linearizability. It is also possible to construct an explicit scenario that violates linearizability.

3.3 Stacks

These results are analogous to those for FIFO queues with Pop in place of Deq and Push in place of Enq.

THEOREM 3.10. *If S is a stack with at least two pushers and a distinct popper, then, for any linearizable implementation of S , $|Push(S)| \geq (u/2)$.*

THEOREM 3.11. *There exists a sequentially consistent implementation of stacks with $|Push| = 0$ and $|Pop| = h$.*

4. PERFECT CLOCKS

In the previous section, we have shown a gap between the cost of implementing sequential consistency and the cost of implementing linearizability. This separation hinged on the assumption that clocks are not initially synchronized. In this section we show that this assumption is necessary by considering the case in which processes have perfectly synchronized (*perfect*) clocks

```

Enqp(Q, v):
  ABC-send(Q, v, "enq")
  generate Ackp(Q)

Deqp(Q):
  ABC-send(Q, "deq")

ABC-receive(Q, v, "enq") from q:
  enqueue v on local copy of Q

ABC-receive(Q, "deq") from q:
  val := dequeue local copy of Q
  if p = q then generate Retp(Q, val) endif

```

Fig. 5. Sequentially consistent fast-enqueue algorithm.

and message delay is constant and known.⁹ Another contribution of these results is that our lower bounds for this model also hold a fortiori in more realistic models. Perfect clocks are modeled by letting $C_p(t) = t$ for all p and t . The constant message delay is modeled by letting $u = 0$; d is known and can be used by the MCS.

For each of the three object types, we first prove lower bounds on the worst-case response time for sequentially consistent implementations. Since sequential consistency is a weaker condition than linearizability, these bounds also hold for linearizable implementations. Then we present algorithms that achieve linearizability and, hence, sequential consistency, with worst-case response times matching the lower bounds. Section 4.1 considers read/write objects, Section 4.2 considers FIFO queues, and Section 4.3 considers stacks.

4.1 Read / Write Objects

We show in Section 4.1.1 that for sequential consistency the sum of the worst-case response times of read and write operations is at least d , even in this strong model. This is a formalization of a result of Lipton and Sandberg [1988, Theorem 1], making precise the timing assumptions made on the system. We then show in Section 4.1.2 that the lower bound is tight for this model by describing two linearizable algorithms that match the lower bound exactly: In the first algorithm, reads are performed instantaneously, while the worst-case response time for a write is d . In the second algorithm, writes are performed instantaneously, while the worst-case response time for a read is d .

⁹The assumptions that processes have perfect clocks and that message delays are constant (and known) are equivalent. If one assumes that clocks are not necessarily synchronized perfectly (but run at the rate of real time) and that the message delay is constant and known, then a simple algorithm suffices to synchronize the clocks perfectly. If one assumes that clocks are perfectly synchronized and that there is a known upper bound d on message delays, then constant message delays can be easily simulated by time-stamping each message with the clock time of the sender and having each recipient delay any message that arrives with a delay smaller than d until the delay is exactly d .

4.1.1 Lower Bounds for Sequential Consistency

THEOREM 4.1 [Lipton and Sandberg 1988]. *For any MCS that is a sequentially consistent implementation of two read/write objects X and Y , $|Write| + |Read| \geq d$.*

PROOF. Let p and q be two processes that access X and Y . Assume by contradiction that there exists a sequentially consistent implementation of X and Y for which both $|Write(X)| + |Read(Y)| < d$ and $|Write(Y)| + |Read(X)| < d$. Without loss of generality, assume that 0 is the initial value of both X and Y .

By the specification of Y , there is some admissible execution α_1 such that $ops(\alpha_1)$ is

$$Write_p(X, 1) Ack_p(X) Read_p(Y) Ret_p(Y, 0),$$

$Write_p(X, 1)$ occurs at real time 0, and $Read_p(Y)$ occurs immediately after $Ack_p(X)$. By assumption, the real time at the end of α_1 is less than d . Thus, no message is received at any node during α_1 .

By the specification of X , there is some admissible execution α_2 such that $ops(\alpha_2)$ is

$$Write_q(Y, 1) Ack_q(Y) Read_q(X) Ret_q(X, 0),$$

$Write_q(Y, 1)$ occurs at real time 0, and $Read_q(X)$ occurs immediately after $Ack_q(Y)$. By assumption, the real time at the end of α_2 is less than d . Thus, no message is received at any node during α_2 .

Since no message is ever received in α_1 and α_2 , the execution α obtained from α_1 , by replacing q 's history with q 's history in α_2 , is admissible. Then $ops(\alpha)$ consists of the operations $[Write_p(X, 1), Ack_p(X)]$ followed by $[Read_p(Y), Ret_p(Y, 0)]$, and $[Write_q(Y, 1), Ack_q(Y)]$ followed by $[Read_q(X), Ret_q(X, 0)]$.

By assumption, α is sequentially consistent. Thus, there is a legal operation sequence τ consisting of the operations $[Write_p(X, 1), Ack_p(X)]$ followed by $[Read_p(Y), Ret_p(Y, 0)]$, and $[Write_q(Y, 1), Ack_q(Y)]$ followed by $[Read_q(X), Ret_q(X, 0)]$. Since τ is a sequence of operations, either the read of X follows the write of X , or the read of Y follows the write of Y . But each possibility violates the serial specification of either X or Y , contradicting τ being legal. \square

4.1.2 Upper Bounds for Linearizability. In this section we show that the trade-off suggested by Theorem 4.1 is inherent and that a sequentially consistent implementation may choose which operation to slow down. More precisely, we present an algorithm in which a read operation is instantaneous (local), while a write operation returns within time d ; we also present an algorithm in which the roles are reversed. These algorithms actually ensure the stronger condition of linearizability.

The algorithm for fast reads and slow writes works as follows: Each process keeps a copy of all objects in its local memory. When a $Read_p(X)$ occurs, p reads the value v of X in its local memory and immediately does a $Ret_p(X, v)$.

When a $\text{Write}_p(X, v)$ occurs, p sends “write(X, v)” messages to all other processes. Then p waits d time, after which it changes the value of X to v in its local memory and does an $\text{Ack}_p(X)$. Whenever a process receives a “write(X, v)” message, it changes the value of X to v in its local memory. (If it receives several at the same time, it “breaks ties” using sender IDs; that is, it writes the value in the message from the process with the largest ID and ignores the rest of the messages.)

THEOREM 4.2. *There exists a linearizable implementation of read/write objects with $|\text{Read}| = 0$ and $|\text{Write}| = d$.*

PROOF. Consider the algorithm just described. Clearly, the time for every read is 0, and the time for every write is d .

Let σ be an admissible execution of this algorithm. For each operation in σ , say that it *occurs* at the real time when its response happens. Let τ be the sequence of operations in σ ordered by time of occurrence, breaking ties with process IDs. Clearly, $\text{ops}(\sigma)|_p$ is equal to $\tau|_p$ for all p , and the order of nonoverlapping operations is preserved.

It remains to show that τ is legal, that is, for every object X , $\tau|X$ is in the serial specification of X . Since X is a read/write object, we must show that every Read returns the value written by the latest preceding Write (and if there is no such Write, then it returns the initial value).

Pick any X , and consider $\tau|X = op_1 op_2 \dots$. Suppose op_i is $[\text{Read}_p(X), \text{Ret}_p(X, v)]$ and op_i occurs at time t in σ .

Case 1. No Write precedes op_i in τ . By the definition of τ , no Write is Acked before op_i starts. Since the Ack for a Write happens at the same time that every process updates its local copy of X , the Read reads the initial value for X and returns that value.

Case 2. Some $\text{Write}_p(X, v)$ is the latest Write preceding op_i in τ . By the definition of τ , this Write is Acked before op_i starts, but no other Write is Acked before op_i starts. Since the Ack for a Write happens at the same time that every process updates its local copy of X , the Read reads v for the value of X and returns that value. \square

The algorithm for slow reads and fast writes is similar to the previous one. Each process keeps a copy of all objects in its local memory. When a $\text{Read}_p(X)$ occurs, p waits d time, after which it reads the value v of X in its local memory and immediately does a $\text{Ret}_p(X, v)$. When a $\text{Write}_p(X, v)$ occurs, p sends “write(X, v)” messages to all other processes (including a dummy message to itself, which is delayed d time) and does an Ack immediately. Whenever a process receives a “write(X, v)” message, it changes the value of X to v in its local memory. Ties are resolved as in the previous algorithm.

THEOREM 4.3. *There exists a linearizable implementation of read/write objects with $|\text{Read}| = d$ and $|\text{Write}| = 0$.*

PROOF. Consider the algorithm just described. Clearly, the time for every read is d , and the time for every write is 0.

Let σ be an admissible execution of this algorithm. For each operation in σ , say that it *occurs* at the real time when its call happens. Let τ be the sequence of operations in σ ordered by time of occurrence, breaking ties with process IDs. Clearly, $ops(\sigma)|_p$ is equal to $\tau|_p$ for all p , and the order of nonoverlapping operations is preserved.

It remains to show that τ is legal, that is, for every object X , $\tau|X$ is in the serial specification of X . Since X is a read write object, we must show that every Read returns the value written by the latest preceding Write (and if there is no such Write, then it returns the initial value).

Pick any X , and consider $\tau|X = op_1op_2\dots$. Suppose op_i is $[Read_p(X), Ret_p(X, v)]$ and op_i occurs at time t in σ .

Case 1. No Write precedes op_i in t . By the definition of τ , no Write starts before op_i starts. Since the local changes occur d time after the Write starts and the Read reads the local memory d time after the Read starts, it reads the local memory before any change is made to it. Thus, the Read returns the initial value.

Case 2. Some $Write_p(X, v)$ is the latest Write preceding op_i in τ . Essentially, the same argument as in Case 1 works. \square

4.2 FIFO Queues

We show in Section 4.2.1 that for sequential consistency the worst-case response time of a dequeue operation is at least d , even when clocks are perfectly synchronized and message delays are constant. We then show in Section 4.2.2 that this lower bound is tight for this model by describing a linearizable algorithm that matches the lower bound exactly: Enqueues are performed instantaneously, while dequeues take time d .

4.2.1 Lower Bound for Sequential Consistency.

THEOREM 4.4. *For any sequentially consistent implementation of a FIFO queue Q , $|Deq(Q)| \geq d$.*

PROOF. Let p and q be two processes that access Q . Assume by contradiction that there exists a sequentially consistent implementation of Q for which $|Deq(Q)| < d$. Let $T = |Deq(Q)|$. By definition, the queue Q is initially empty.¹⁰ By the specification of Q , there is some admissible execution α'_1 such that $ops(\alpha'_1)$ is

$$\begin{aligned} & \text{Enq}_q(Q, 1) \text{Ack}_q(Q) \text{Deq}_p(Q) \text{Ret}_p(Q, v_1) \dots \\ & \text{Deq}_p(Q) \text{Ret}_p(Q, v_i) \dots, \end{aligned}$$

$\text{Enq}_q(Q, 1)$ occurs at real time 0, and $\text{Ack}_q(Q)$ occurs at time t ; the first $\text{Deq}_p(Q)$ occurs at time t , while the j th $\text{Deq}_p(Q)$ occurs at time $t + (j - 1)T$

¹⁰If we allow queues to be initially nonempty, the proof of the lower bound becomes much simpler; we leave the details to the interested reader.

(see Figure 6a). Consider now the infinite sequence v_1, \dots, v_i, \dots . It is possible that many of them are \perp ; however, since only a finite number of Deq operations can be serialized before the Enq operation, we have the following lemma:

LEMMA 4.5. *There exists some i such that $v_i \neq \perp$.*

Fix this particular i , and note that $v_i = 1$ and, for all j , $1 \leq j < i$, $v_j = \perp$. Let α_1 be α'_1 truncated after the i th Deq operation by p . More precisely, $ops(\alpha_1)$ is

$$\begin{aligned} & \text{Enq}_q(Q, 1) \text{Ack}_q(Q) \text{Deq}_p(Q) \text{Ret}_p(Q, \perp) \cdots \\ & \text{Deq}_p(Q) \text{Ret}_p(Q, \perp) \text{Deq}_p(Q) \text{Ret}_p(Q, 1). \end{aligned}$$

$\text{Enq}_q(Q, 1)$ occurs at real time 0, and $\text{Ack}_q(Q)$ occurs at time t ; the first $\text{Deq}_p(Q)$ occurs at time t , while the i th $\text{Deq}_p(Q)$ occurs at time $t + (i - 1)T$ (see Figure 6b). It is clear that the v_j 's are exactly as in α'_1 . By assumption, the real time at the end of α is less than $t + (i - 1)T + d$. Thus, no message sent after $t + (i - 1)T$ is received during α_1 .

We now consider the execution where the i th (and last) dequeue by p is replaced with a dequeue by q . More precisely, by the specification of Q , there is some admissible execution α_2 such that $ops(\alpha_2)$ is

$$\begin{aligned} & \text{Enq}_q(Q, 1) \text{Ack}_q(Q) \text{Deq}_p(Q) \text{Ret}_p(Q, \perp) \cdots \\ & \text{Deq}_p(Q) \text{Ret}_p(Q, \perp) \text{Deq}_q(Q) \text{Ret}_q(Q, u), \end{aligned}$$

$\text{Enq}_q(Q, 1)$ occurs at real time 0, and $\text{Ack}_q(Q)$ occurs at time t ; the first $\text{Deq}_p(Q)$ occurs at time t , while the $(i - 1)$ st $\text{Deq}_p(Q)$ occurs at time $t + (i - 2)T$, and Deq_q occurs at time $t + (i - 1)T$ (see Figure 6c). Since α_2 is sequentially consistent, it follows that $u = 1$. By assumption, the real time at the end of α_2 is less than $t + (i - 1)T + d$. Thus, no message sent after $t + (i - 1)T$ is received during α_2 .

Consider now an execution α obtained from α_1 by replacing q 's history with q 's history in α_2 . No message sent after time $t + (i - 1)T$ is ever received in α_1 or α_2 , and α_1 and α_2 are identical until time $t + (i - 1)T$. This implies that α is admissible. Then $ops(\alpha)$ is

$$\begin{aligned} & \text{Enq}_q(Q, 1) \text{Ack}_q(Q) \text{Deq}_p(Q) \text{Ret}_p(Q, \perp) \cdots \\ & \text{Deq}_p(Q) \text{Ret}_p(Q, 1) \text{Deq}_q(Q), \text{Ret}_q(Q, 1) \end{aligned}$$

(see Figure 6d). By assumption, α is sequentially consistent. Thus, there is a legal sequence τ , which is a permutation of the above operations. However, in τ the element "1" is enqueued once, but dequeued twice, a contradiction. \square

4.2.2 Upper Bound for Linearizability. In this section we show that the lower bound given in Theorem 4.4 is tight for the model with perfect clocks. Specifically, we present an algorithm in which an enqueue operation returns instantaneously, while a dequeue operation returns within time d . The algorithm ensures the stronger condition of linearizability.

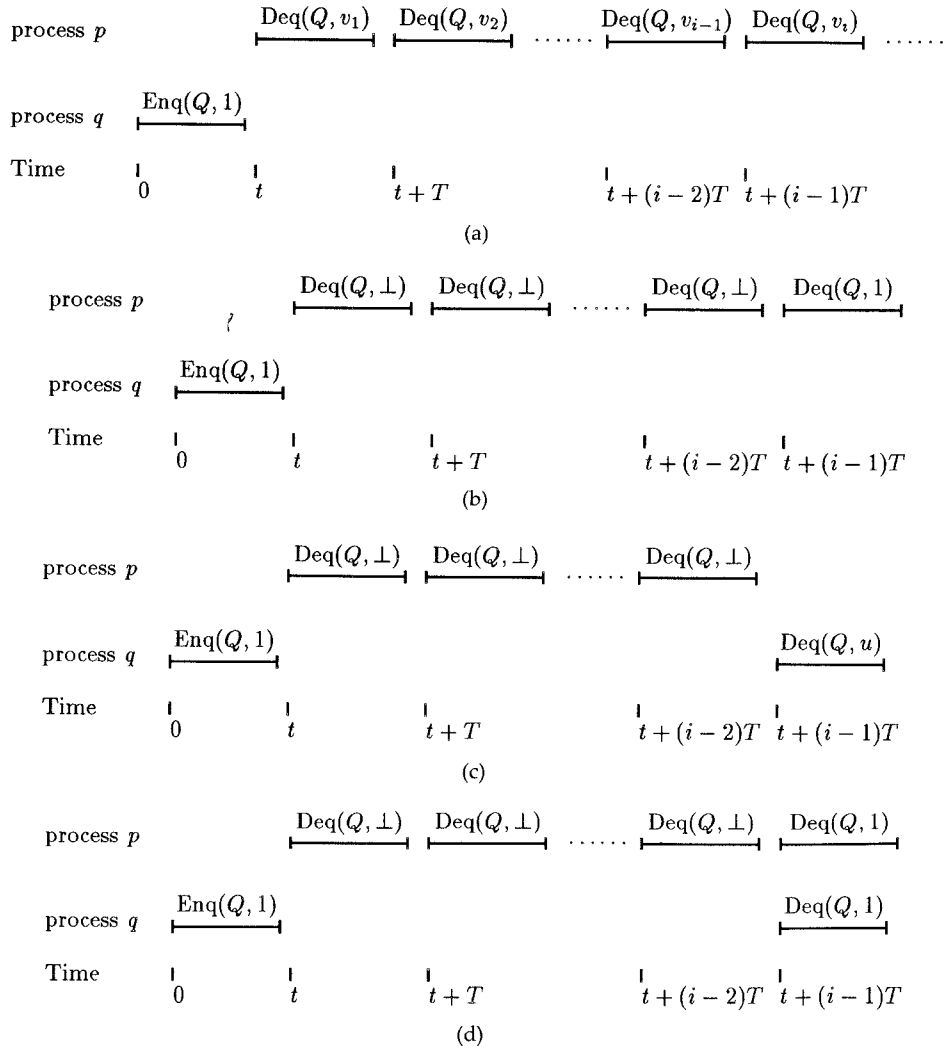


Fig. 6. Executions used in the proof of Theorem 4.4 (a) Execution α'_1 . (b) Execution α_1 . (c) Execution α_2 . (d) Execution α .

The algorithm works as follows: Each process keeps a copy of all queues in its local memory. When an $\text{Enq}_p(Q, v)$ occurs, p sends “enqueue(Q, v)” messages to all other processes (including a message to itself which is delayed d time) and does an Ack immediately. When a $\text{Deq}_p(Q)$ occurs, p sends “dequeue(Q)” messages to all other processes (including a message to itself which is delayed d time). After waiting d time, p handles its own message and does a $\text{Ret}_p(Q, v)$. Whenever a process receives an “enqueue(Q, v)” or “dequeue(Q)” message, it makes the appropriate update to the copy of Q in its local memory. (If it receives several messages at the same time, it “breaks

ties” using sender IDs; that is, it handles them by increasing order of process IDs.)

THEOREM 4.6. *There exists a linearizable implementation of FIFO queues with $|Enq| = 0$ and $|Deq| = d$.*

In the proof we serialize each operation to occur d time after it is called. Since all processes update their local copies at these serialization times, the claim follows:

PROOF. Consider the algorithm just described. Clearly, $|Enq| = 0$, and $|Deq| = d$.

Let σ be an admissible execution of this algorithm. For each operation in σ , say that it *occurs* at time d after the real time when its call happens. Let τ be the sequence of operations in σ ordered by time of occurrence, breaking ties with process IDs. Clearly, $ops(\sigma)|_p$ is equal to $\tau|_p$ for all p , and the order of nonoverlapping operations is preserved.

It remains to show that τ is legal; that is, for every object Q , $\tau|_Q$ is in the serial specification of Q . Pick any Q , and consider $\tau|_Q = op_1 op_2 \dots$. Suppose op_i is $[Deq_p(Q), Ret_p(Q, v)]$. Because the message delay is fixed, updates at p to the local copy of Q occur in the same order as in τ , and the claim follows. \square

4.3 Stacks

The results for stacks are analogous to those for FIFO queues, with Pop playing the role of Deq and Push the role of Enq.

THEOREM 4.7. *For any sequentially consistent implementation of a stack S , $|Pop(S)| \geq d$.*

THEOREM 4.8. *There exists a linearizable implementation of stacks with $|Push| = 0$ and $|Pop| = d$.*

5. CONCLUSIONS AND FURTHER RESEARCH

We have presented a quantitative comparison of the data access time for two well-known consistency conditions for concurrently accessed shared data: sequential consistency and linearizability. Our results indicate that supporting sequential consistency can be more cost effective than supporting linearizability, for certain object types and under certain timing assumptions. Our results also show that a very precise definition of the guarantees provided is important, since seemingly minor differences in the definitions result in significant differences in the inherent efficiency of implementing them. Since our lower bounds are proved in a very strong model, they clearly hold for more practical systems. We believe our algorithms can be adapted to work in more realistic systems.

Our work is closely related to the design of cache consistency schemes that guarantee sequential consistency [Brantley et al. 1985; Censier and Feautrier 1978; Collier 1984; Dubois et al. 1986b; Lamport 1979]. Our implementations use ideas similar to those previously used in cache coherence protocols (cf.,

Archibald and Baer [1986]). In a sense, we have snooping protocols (without a bus) with a *write broadcast* policy (cf., Hennessey and Patterson [1990, pp. 467–469]). We believe our ideas can be modified to accommodate a *write invalidate* policy; however, this will slow down the reads.

Our results can be extended to obtain bounds on the response time of implementing other objects, for example, Test & Set registers, under sequential consistency and linearizability. Further work in this direction is currently under way [Friedman 1993; Kosa 1994].

The modular usage of atomic broadcast in our implementations of sequential consistency admits several extensions. For example, a bus provides an easy mechanism for atomic broadcast by enforcing a global ordering on all messages delivered to the processes. Afek et al. [1993] present a sequentially consistent implementation of read/write objects, for systems where processes communicate via a bus. It might be possible to improve and simplify the correctness proof of the algorithm in Afek et al. [1993] using this observation. Also, atomic broadcast algorithms can be made *fault-tolerant*. This can help in the design of MCSs that can sustain failures of some of the processes. In general, the issue of fault-tolerance is rarely addressed in the current research on memory consistency. As multiprocessors scale up and the probability of failure increase, this will become an important concern.

Following the original appearance of our results [Attiya and Welch 1991], Mavronicolas and Roth [1992] have shown that the trade-off $|W| + |R| = d$ for sequential consistency is indeed continuous and that there are algorithms that achieve all intermediate values. Furthermore, they have shown that the $\Omega(u)$ bounds we prove for linearizability in the approximate clocks model are tight, by extending our algorithms for the perfect-clocks model. It is known that u depends on how closely the clocks in the system are synchronized. Since closely synchronized clocks admit more efficient implementations of linearizability, it may be worthwhile to provide such clocks.

Recently, several nonglobal conditions that are weaker than sequential consistency have been suggested, for example, weak ordering [Adve and Hill 1990b; Bisiani et al. 1989; and Dubois et al. 1986a]; release consistency [Gharachorloo et al. 1990; Gibbons et al. 1991]; pipelined memory [Lipton and Sandberg 1988]; slow memory [Hutto and Ahamad 1989]; causal memory [Ahamad et al. 1990]; loosely coherent memory [Bennett et al. 1990]; and the definitions in Collier [1984] and Ramachandran et al. [1989]. It would be interesting to investigate the inherent efficiency of supporting these consistency guarantees. In order to do so, crisp and precise definitions of these conditions are needed. Results in this direction appear in Attiya and Friedman [1992] and in Attiya et al. [1993].

The cost measure we have chosen to analyze is response time, but it is clear that efficiency, in general, and response time, in particular, are not the only criteria for evaluating consistency guarantees. Other important quantitative measures are amount of local processing and level of message traffic. It has been suggested that it may be possible to implement linearizability *more* cheaply than sequential consistency with regard to these measures, since linearizability is a local property and sequential consistency is not. More

qualitative properties such as the ease of designing, verifying, programming, and debugging algorithms using such shared memories are also very important.

Our formal model ignores several important practical issues, for example, limitations on the size of local memory storage, network topology, clock drift, and “hot-spots.” It will be interesting to understand how these issues influence the bounds.

As multiprocessor systems become larger, distributed implementations of shared virtual memory are becoming more common, since truly shared memories, or even buses, cannot be used in systems with a large number of processors. Such implementations and their evaluation relate issues concerning multiprocessor architecture, programming language design, software engineering, and the theory of concurrent systems. (For instance, our work makes use of shifting and atomic broadcast techniques from the theoretical and practical distributed computing literature.) We hope our work contributes toward a more solid ground for this interaction.

APPENDIX

Atomic Broadcast

The atomic broadcast algorithm employed by our algorithms is based on assigning time stamps to messages. Each process maintains a local time stamp (counter) and a vector with (conservative) estimates of the time stamps of all other processes. A process keeps a time stamp bigger than or equal to the time stamps of all the other processes (according to its estimates). Upon a request to broadcast a message, the message is tagged with the requestor’s current time stamp. Each process maintains a set of messages that are waiting to be delivered. A message with time stamp x is delivered only when the process is certain that all other messages with a time stamp $\leq x$ have arrived at it. This is done by waiting to learn that all processes have increased their time stamps to be at least $x + 1$.¹¹ Once it learns that all processes have increased their time stamps beyond x , the process handles all pending messages with time stamps less than or equal to x , in order, breaking ties using process IDs.

More precisely, to broadcast a message m , p sends a message (t_p, m) to all processes (including itself), where t_p is p ’s current time stamp. It then increases its own time stamp by one and returns. When a process q receives a message with time stamp t_p from p , it saves it in a list of pending messages, sorted by time stamp and process ID. It then increases its time stamp to be at least as large as $t_p + 1$ and sends a time stamp increase message “timestamp(t_q, q)”.

When a process receives a time stamp increase message, it update the time stamp entry for the sender and checks to see if there are any pending

¹¹For simplicity, the algorithm presented here assumes FIFO channels. This assumption can be removed if sequence numbers are employed.

messages whose time stamp is strictly less than all processes' time stamps (saved in its local vector). These messages are delivered in increasing time stamp order, breaking ties using process IDs.

The algorithm uses the following data types:

time stamp = integer,
message = record with fields.
 mess: string (message to be delivered),
 ts: time stamp (assigned by initiator),
 ID: process ID (ID of initiator).

Each process knows n , the total number of processes.

The state of each process consists of the following components:

ts: array[1... n] of integer, all initially 0
 (estimate (from below) time stamps of all processes),
pending: set of message, initially empty
 (set of messages waiting to be delivered).

The transition function of process p appears in Figure 7.

To show that this algorithm implements atomic broadcast, we must show, for any admissible execution, that messages are delivered in the same order to all processes. The ordering of messages is done by time stamps (breaking ties with process IDs). The resulting sequence respects the order at each process by construction and because of the way time stamps are assigned.

More formally, fix some admissible execution σ of the algorithm. The next lemma follows immediately from the code:

LEMMA A.1. *Let p be any process. Then every message broadcast by p in σ is given a unique time stamp in increasing order.*

This immediately implies Lemma A.2:

LEMMA A.2. *The time stamps assigned to messages in σ , together with process IDs, form a total order.*

This total order is called *time stamp order*.

LEMMA A.3. *Let p be any process. Then all messages are delivered to p in σ in time stamp order.*

PROOF. Let (t_1, q_1) be the time stamp of the message m_1 , and let (t_2, q_2) be the time stamp of the message m_2 . Suppose, by contradiction, that $(t_1, q_1) < (t_2, q_2)$, but m_2 was delivered to p before m_1 .

When m_2 is delivered to p , it cannot yet have the message m_1 in pending because otherwise it would deliver it before m_2 . By the code, in order to deliver m_2 , it must be that $ts_p[q_1] > t_2$. But then p must have received a time-stamp message from q_1 with a time stamp $t \geq t_2 + 1$. Since $(t_1, q_1) < (t_2, q_2)$, it must be that $t_1 \leq t_2$, and hence, $t > t_1$. By the code, the message

```

ABC-sendp(m):
  send (ts[p],m) to all processes
  ts[p] := ts[p] + 1

receive (t,m) from q:
  add (m,t,q) to pending
  if t + 1 > ts[p] then
    ts[p] := t + 1
    send timestamp(ts[p]) to all processes
  endif

receive timestamp(t) from q:
  ts[q] := t
  repeat
    let E be element with smallest (ts,id) pair in pending
    if for some q, ts[q] ≤ E.ts then exit
    deliver E.m { this is the ABC-receive }
    remove E from pending
  endrepeat

```

Fig. 7. Atomic broadcast algorithm.

m_1 was sent before the time-stamp message. But then the FIFO property of the communication system implies that p has already received m_1 —a contradiction. \square

The next lemma guarantees that each message is delivered within time $2d$ from the initiation of the operation.

LEMMA A.4. *If process p broadcasts a message m , then m is delivered at each process within time at most $2d$ in σ .*

PROOF. Assume p broadcasts m at time T , with time stamp x . By time $T + d$, all processes will get the message (x, m) and will set their time stamps to be at least $x + 1$, sending a time-stamp increase message to all other processes, if necessary. Thus, by time $T + 2d$, all processes will have in their time-stamp vectors values that are strictly larger than x and will deliver m . \square .

Lemmas A.3 and A.4 prove the following theorem:

THEOREM A.5. *The algorithm in Figure 7 is an atomic broadcast algorithm with $h = 2d$.*

ACKNOWLEDGMENTS

The authors thank Sarita Adve, Roy Friedman, Mark Hill, and Rick Zucker for helpful comments on an earlier version of this paper. We especially thank Martha Kosa for a careful reading. The comments of the anonymous referees helped us improve the presentation.

REFERENCES

- ADVE, S., AND HILL, M. 1990a. Implementing sequential consistency in cache-based systems. In *Proceedings of the International Conference on Parallel Processing*. Pennsylvania State University, University Park, pp. I-47–I-50.
- ADVE, S., AND HILL, M. 1990b. Weak ordering—A new definition. In *Proceedings of the 17th International Symposium on Computer Architecture* (Seattle, Wash., May 21–31). IEEE, Los Angeles, pp. 2–14.
- AFEK, Y., BROWN, G., AND MERRITT, M. 1993. Lazy caching. *ACM Trans. Program. Lang. Syst.* 15, 1 (Jan.), 182–205.
- AHAMAD, M., HUTTO, P. AND JOHN, R. 1990. Implementing and programming causal distributed shared memory, Tech. Rep. GIT-CC-90-49, Georgia Inst. of Technology, Atlanta, Dec.
- ARCHIBALD, J., AND BAER, J.-L. 1986. Cache coherence protocols: Evaluation using a multiprocessor simulation model, *ACM Trans. Comput. Syst.* 4, 4 (Nov.), 273–298.
- ATTIYA, H. 1991. Implementing FIFO queues and stacks. In *Proceedings of the 5th International Workshop on Distributed Algorithms*. Lecture Notes in Computer Science, vol. 579, Springer-Verlag, New York, pp. 80–94. (Also, Tech. Rep. 672, Dept. of Computer Science, The Technion, Haifa, May).
- ATTIYA, H., AND FRIEDMAN, R. 1992. A consistency condition for high-performance multiprocessors. In *Proceedings of the 24th ACM Symposium on Theory of Computing* (Victoria, B.C., May 4–6). ACM, New York, pp. 679–690. (Also, Tech. Rep. 719, Dept. of Computer Science, The Technion, Haifa, June 1992).
- ATTIYA, H., AND WELCH, J. L. 1991. Sequential consistency versus linearizability. In *Proceedings of the 3rd ACM Symposium on Parallel Algorithms and Architectures* (Hilton Head, S.C., July 21–24). ACM, New York, pp. 304–315.
- ATTIYA, H., CHAUDHURI, S., FRIEDMAN, R., AND WELCH, J. L. 1993. Non-sequential consistency conditions for shared memory. In *Proceedings of the 5th ACM Symposium on Parallel Algorithms and Architectures* (Velen, Germany, June 30–July 2). ACM, New York, pp. 241–250.
- BENNETT, J., CARTER, J., AND ZWAENEPOEL, W. 1990. Munin: Distributed shared memory based on type-specific memory coherence. In *Proceedings of the 2nd ACM Symposium on Principles and Practice of Parallel Processing* (Seattle, Wash., Mar. 14–16). ACM, New York, pp. 168–176.
- BERNSTEIN, P., HADZILACOS, V., AND GOODMAN, H. 1987. *Concurrency Control and Recovery in Database Systems*. Addison-Wesley, Reading, Mass.
- BIRMAN, K., AND JOSEPH, T. 1987. Reliable communication in the presence of failures. *ACM Trans. Comput. Syst.* 5, 1 (Feb.), 47–76.
- BISIANI, R., NOWATZYK, A., AND RAVISHANKAR, M. 1989. Coherent shared memory on a distributed memory machine. In *Proceedings of the International Conference on Parallel Processing*. Pennsylvania State University, University Park, pp. I-133–I-141.
- BRANTLEY, W., MCAULIFFE, K., AND WEISS, J. 1985. RP3 processor-memory element. In *Proceedings of the International Conference on Parallel Processing* (Aug. 20–23). Pennsylvania State University, University Park, pp. 782–789.
- CENSIER, L. M., AND FEAUTRIER, P. 1978. A new solution to coherence problems in multicache systems. *IEEE Trans. Comput. C-27*, 12 (Dec.), 1112–1118.
- CHANG, J., AND MAXEMCHUK, N. F. 1984. Reliable broadcast protocols. *ACM Trans. Comput. Syst.* 2, 3 (Aug.), 251–273.
- COLLIER, W. W. 1984. Architectures for systems of parallel processes. Tech. Rep. 00.3253, IBM, Poughkeepsie, N.Y., Jan.
- DUBOIS, M., AND SCHEURICH, C. 1990. Memory access dependencies in shared-memory multiprocessors. *IEEE Trans. Softw. Eng.* 16, 6 (June), 660–673.
- DUBOIS, M., SCHEURICH, C., AND BRIGGS, F. A. 1986a. Memory access buffering in multiprocessors. In *Proceedings of the 13th International Symposium on Computer Architecture* (June), pp. 434–442.
- DUBOIS, M., SCHEURICH, C., AND BRIGGS, F. A. 1986b. Synchronization, coherence and event ordering in multiprocessors. *Computer* 21, 2 (June), 9–21.
- FRIEDMAN, R. 1993. Implementing hybrid consistency with high-level synchronization opera-

- tions, *Proceedings of the 12th ACM Symposium on Principles of Distributed Computing* (Ithaca, N.Y., Aug. 15–18). ACM, New York, pp. 229–240.
- GARCIA-MOLINA, H., AND SPAUSTER, A. 1989. Message ordering in a multicast environment. In *Proceedings of the International Conference on Distributed Computing Systems* (Newport Beach, Calif., June 5–9). IEEE, Los Angeles, pp. 354–361.
- GHARACHORLOO, K., LENOSKI, D., LAUDON, J., GIBBONS, P., GUPTA, A., AND HENNESSEY, J. 1990. Memory consistency and event ordering in scalable shared-memory multiprocessors. In *Proceedings of the 17th International Symposium on Computer Architecture* (Seattle, Wash., May 28–31). IEEE, Los Angeles, pp. 15–26.
- GIBBONS, P., MERRITT, M., AND GHARACHORLOO, K. 1991. Proving sequential consistency of high-performance shared memories. In *Proceedings of the 3rd ACM Symposium on Parallel Algorithms and Architectures* (Hilton Head, S.C., July 21–24). ACM, New York, pp. 292–303.
- GOTTLIEB, A., GRISHMAN, R., KRUSKAL, C. K., MCAULIFFE, K. P., RUDOLPH, L., AND SNIR, M. 1983. The NYU ultracomputer—Designing a MIMD shared-memory parallel computer. *IEEE Trans. Comput. C-32* 3 (Feb.) 175–189.
- HENNESSEY, J., AND PATTERSON, D. 1990. *Computer Architecture: A Quantitative Approach*. Morgan Kaufmann, San Mateo, Calif.
- HERLIHY, M. 1988. Wait-free implementations of concurrent objects. In *Proceedings of the ACM Symposium on Principles of Distributed Computing* (Toronto, Ont., Aug. 15–17). ACM, New York, pp. 276–290.
- HERLIHY, M., AND WING, J. 1990. Linearizability: A correctness condition for concurrent objects. *ACM Trans. Program. Lang. Syst.* 12, 3 (July), 463–492.
- HUTTO, P., AND AHAMAD, M. 1989. Slow memory: Weakening consistency to enhance concurrency in distributed shared memories. Tech. Rep. GIT-ICS-89/39, Georgia Inst. of Technology, Atlanta, Oct.
- KOSA, M. 1994. Consistency conditions for concurrent shared objects: Upper and lower bounds. Ph.D. dissertation, Dept. of Computer Science, Univ. of North Carolina, Chapel Hill, Feb.
- LAMPORT, L. 1979. How to make a multiprocessor computer that correctly executes multiprocess programs. *IEEE Trans. Comput. C-28*, 9 (Sept.), 690–691.
- LAMPORT, L. 1986. On interprocess communication. Parts I and II. *Distrib. Comput.* 1, 2, 77–101.
- LIPTON, R., AND SANDBERG, J. 1988. PRAM. A scalable shared memory. Tech. Rep. CS-TR-180-88, Dept. of Computer Science, Princeton Univ., Princeton, N.J., Sept.
- LUNDELIUS, J., AND LYNCH, N. 1984. An upper and lower bound for clock synchronization. *Inf. Control* 62, 2–3 (Aug/Sept.), 190–204.
- MAVRONICOLAS, M., AND ROTH, D. 1992. Efficient, strong consistent implementations of shared memory. In *Proceedings of the 6th International Workshop on Distributed Algorithms*. Lecture Notes in Computer Science, vol. 647, Springer-Verlag, New York, pp. 346–361.
- MIN, S., AND BAER, J. 1989. A timestamp-based cache coherence scheme. In *Proceedings of the International Conference on Parallel Processing*. Pennsylvania State University, University Park, pp. I-23–I-32.
- MISRA, J. 1986. Axioms for memory access in asynchronous hardware systems. *ACM Trans. Program. Lang. Syst.* 8, 1 (Jan.), 142–153.
- PAPADIMITRIOU, C. 1986. *The Theory of Concurrency Control*. Computer Science Press, Rockville, Md.
- RAMACHANDRAN, U., AHAMAD, J., AND KHALIDI, M. Y. 1989. Coherence of distributed shared memory: Unifying synchronization and data transfer. In *Proceedings of the International Conference on Parallel Processing*. Pennsylvania State University, University Park, pp. II-160–II-169.
- SCHEURICH, C., AND DUBOIS, M. 1987. Correct memory operation of cache-based multiprocessors. In *Proceedings of the 14th International Symposium on Computer Architecture* (Pittsburgh, Pa., June 2–5). IEEE, Los Angeles, pp. 234–243.

Received August 1991; revised March 1993; accepted May 1993