

On the Space Complexity of Randomized Synchronization

FAITH FICH

University of Toronto, Toronto, Ont., Canada

MAURICE HERLIHY

Brown University, Providence, Rhode Island

AND

NIR SHAVIT

Tel-Aviv University, Tel-Aviv, Israel

Abstract. The “wait-free hierarchy” provides a classification of multiprocessor synchronization primitives based on the values of n for which there are deterministic wait-free implementations of n -process consensus using instances of these objects and *read-write* registers. In a randomized wait-free setting, this classification is degenerate, since n -process consensus can be solved using only $O(n)$ *read-write* registers.

In this paper, we propose a classification of synchronization primitives based on the *space complexity* of randomized solutions to n -process consensus. A *historyless* object, such as a *read-write* register, a *swap* register, or a *test&set* register, is an object whose state depends only on the last nontrivial operation that was applied to it. We show that, using *historyless* objects, $\Omega(\sqrt{n})$ object instances are necessary to solve n -process consensus. This lower bound holds even if the objects have unbounded size and the termination requirement is *nondeterministic solo termination*, a property strictly weaker than randomized wait-freedom.

This work was supported by the Natural Science and Engineering Research Council of Canada grant A1976 and the Information Technology Research Centre of Ontario, ONR grant N00014-91-J-1046, NSF grant 8915206-CCR and 9520298-CCR, and DARPA grants N00014-92-J-4033 and N00014-91-J-1698.

A preliminary version of this paper appears in the *Proceedings of 12th Annual ACM Symposium on Principles of Distributed Computing* (Ithaca, N.Y., Aug. 15–18). ACM, New York, pp. 241–249.

Part of this work was performed while Faith Fich and Nir Shavit were visiting MIT and Maurice Herlihy was at DEC Cambridge Research Laboratory.

Authors’ addresses: F. Fich, Department of Computer Science, University of Toronto, Toronto, Ontario, Canada M5S 3G4, e-mail: fich@cs.toronto.edu; M. Herlihy, Computer Science Department, Brown University, Box 1910, Providence, RI 02912, e-mail: herlihy@cs.brown.edu; N. Shavit, Computer Science Department Tel-Aviv University, Tel-Aviv 69978, Israel, e-mail: shavir@math.tau.ac.il.

Permission to make digital/hard copy of part or all of this work for personal or classroom use is granted without fee provided that the copies are not made or distributed for profit or commercial advantage, the copyright notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery (ACM), Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists, requires prior specific permission and/or a fee.

© 1998 ACM 0004-5411/98/1100-0843 \$05.00

We then use this result to relate the randomized space complexity of basic multiprocessor synchronization primitives such as *shared counters*, *fetch&add* registers, and *compare&swap* registers. Viewed collectively, our results imply that there is a separation based on space complexity for synchronization primitives in randomized computation, and that this separation differs from that implied by the deterministic “wait-free hierarchy.”

Categories and Subject Descriptors: F.1.2 [Computation by Abstract Devices]: Modes of Computation—*parallelism and concurrency*; F.1.3 [Computation by Abstract Devices]: Complexity Measures and Classes—*complexity hierarchies*

General Terms: Algorithms, Theory

Additional Key Words and Phrases: Consensus, lower bounds, space complexity

1. Introduction

Traditionally, the theory of interprocess synchronization has centered around the notion of *mutual exclusion*: ensuring that only one process at a time is allowed to modify complex shared data objects. As a result of the growing realization that unpredictable delay is an increasingly serious problem in modern multiprocessor architectures, a new class of wait-free algorithms have become the focus of both theoretical [Lynch 1996] and experimental research [Alemany and Felten 1992; Bershad 1993]. An implementation of a concurrent object is *wait-free* if it guarantees that every process will complete an operation within a finite number of its own steps, independent of the level of contention and the execution speeds of the other processes. Wait-free algorithms provide the additional benefit of being highly fault-tolerant, since a process can complete an operation even if all $n - 1$ others fail by halting.

The consensus problem is a computational task that requires each of n asynchronous processes, with possibly different private input values, to decide on one of the inputs as their common output. The “wait-free hierarchy” [Herlihy 1991] classifies concurrent objects and multiprocessor synchronization primitives based on the values of n for which they solve n -process consensus in a wait-free manner. For example, it is impossible to solve n -process consensus using *read-write* registers for $n > 1$ or using *read-write* registers and *swap* registers, for $n > 2$ [Abrahamson 1988; Chor et al. 1987; Dolev et al. 1987; Herlihy 1991b; Loui and Abu-Amara 1987]. It has been shown that this separation does not hold in a randomized setting; that is, even *read-write* registers suffice to solve n -process consensus [Abrahamson 1988; Chor et al. 1987]. This is a rather fortunate outcome, since it opens the possibility of using randomization to implement concurrent objects without resorting to non-resilient mutual exclusion [Aspnes 1990; Aspnes and Herlihy 1990; Aspnes and Waarts 1992; Herlihy 1991a; Saks et al. 1991]. One important application is the software implementation of one synchronization object from another. This allows easy porting of concurrent algorithms among machines with different hardware synchronization support. However, in order to understand the conditions under which such implementations are effective and useful, we must be able to quantify the randomized computational power of existing hardware synchronization primitives.

In this paper, we propose such a quantification by providing a complexity separation among synchronization primitives based on the *space complexity* of randomized wait-free solutions to n -process binary consensus. It is a step towards

a theory that would allow designers and programmers of multiprocessor machines to use mathematical tools to evaluate the power of alternative synchronization primitives and to recognize when certain randomized protocols are inherently inefficient.

Randomized n -process consensus can be solved using $O(n)$ read-write registers of bounded size [Aspnes and Herlihy 1990]. Our main result is a proof that using only *historyless* objects (for example, *read-write* registers of unbounded size, *swap* registers, and *test&set* registers), $\Omega(\sqrt{n})$ instances are necessary to solve randomized n -process binary consensus. We do so by showing a more general result: there is no implementation of consensus satisfying a special property, *nondeterministic solo termination*, from a sufficiently small number of historyless objects. The nondeterministic solo termination property is strictly weaker than randomized wait-freedom.

Our result thus shows that, for n -process consensus, the number of object instances, rather than the sizes of their state spaces, is the important factor. Furthermore, allowing operations such as SWAP or TEST&SET, in addition to READ and WRITE, does not substantially reduce the number of objects necessary to solve consensus.

A variety of lower bounds for randomized algorithms exist in the literature. A lower bound for the Byzantine Agreement problem in a randomized setting was presented by Graham and Yao [1989]. They showed a lower bound on the probability of achieving agreement among 3 processes, one of which might behave maliciously. Their result is derived based on the possibly malicious behavior of a processor.

A randomized space complexity lower bound was presented by Kushilevitz et al. [1993]. They prove lower bounds on the size of (i.e., number of bits in) a read-modify-write register necessary to provide a fair randomized solution to the mutual exclusion problem. Their results relate the size of the register to the amount of fairness in the system and the number of processes accessing the critical section in a mutually exclusive manner. A deterministic space lower bound on the number of bits in a compare&swap register necessary to solve n -process consensus was proved by Afek and Stupp [1993].

Lower bounds on the space complexity of consensus have also been obtained for models in which *objects* as well as processes may fail [Afek et al. 1995; Jayanti et al. 1992].

A powerful time complexity lower bound was recently presented by Aspnes [1997] for the randomized consensus problem when up to t processes may fail by halting. Aspnes studies the total number of coin flips necessary to create a global shared coin, a mathematical structure that he shows must implicitly be found in any protocol solving consensus, unless its complexity is $\Omega(t^2)$. His bound on the shared coin construction implies an interesting $\Omega(t^2/\log^2 t)$ lower bound on the amount of work (total number of operations by all processes) necessary to solve randomized consensus.

Our proof technique is most closely related to the elegant method introduced by Burns and Lynch to prove a lower bound on the number of read/write registers required for a deterministic solution to the mutual-exclusion problem. Though related, the problem they tackle is fundamentally different and in a sense easier than ours. This is because the object they consider (mutual exclusion) is accessed by each process repeatedly, whereas our lower bound applies to the

implementation of general objects and, in particular, consensus (for which each process may perform only a single access).

Because the objects we consider are “single access,” our lower bound proofs required the development of a new proof technique. The key element of this technique is a method of “cutting” and “splicing together” *interruptible executions*, executions that can be broken into pieces between which executions involving other processes can be inserted.

Based on our consensus lower bound, we are able to relate the randomized complexity of basic multiprocessor synchronization primitives such as *counters* [Aspnes and Herlihy 1990; Moran et al. 1992], *fetch&add* registers, and *compare&swap* registers. For example, *swap* registers and *fetch&add* registers are each sufficient to deterministically solve 2-process consensus, but not 3-process consensus. However, a single *fetch&add* register suffices to solve randomized n -process consensus, whereas $\Omega(\sqrt{n})$ *swap* registers are needed. Furthermore, a primitive such as *compare&swap*, which is deterministically “stronger” than *fetch&add*, is essentially equivalent to it under this randomized complexity measure. Our theorems imply that there is a space-complexity-based separation among synchronization primitives in randomized computation and that their relative power differs from what could be expected given the deterministic “wait-free hierarchy.” Our hope is that such separation results will eventually allow hardware designers and programmers to make more informed decisions about which synchronization primitives to use in a randomized setting and how to better emulate one primitive from another.

The structure of our presentation is as follows: In Section 2, we describe our model of computation and how randomization and asynchrony are expressed within it. Section 3 begins by proving a special case of the lower bound, followed by the proof of the general case. Section 4 presents several separation theorems among synchronization primitives based on our main lower bound.

2. Model

Our model of computation consists of a collection of n sequential threads of control called *processes* that communicate by applying operations to shared *objects*.

Objects are data structures in memory. Each object has a *type*, which defines a set of possible *values* and a set of primitive *operations* that provide the only means to manipulate that object. The current value of an object and the operation that is applied to it determine (perhaps nondeterministically) the response to the operation and the (possibly) new value of the object.

For example, a *test&set* register has $\{0, 1\}$ as its set of possible values. Its initial value is 0. The TEST&SET operation responds with the value of the object and then sets the value to 1. A *read-write* register may have any (finite or infinite) set as its set of values. Its operations are READ, which responds with the value of the object, and WRITE(x), for x in the value set, which sets the value of the object to x .

Another example is the *counter* [Aspnes and Herlihy 1990; Moran et al. 1992]. The integers are its set of values. Its operations are INC and DEC, which increment and decrement the counter, respectively, RESET, which sets the value of the counter to 0, and READ, which responds with the value of the counter,

leaving it unchanged. The first three of these operations only respond with a fixed acknowledgement. A *bounded counter* is a counter whose set of possible values is a range of integers and whose operations are performed modulo the size of that range.

Each process has a set of *states*. The operation a process applies and the object to which it is applied depends on the current state of the process. Processes may also have internal operations, such as coin flips. The current state of a process and the result of the operation performed determine the next state of the process. Each such operation is called a *step* of the process. Processes are asynchronous, that is, they can halt or display arbitrary variations in speed. In particular, one process cannot tell whether another has halted or is just running very slowly.

An *execution* is an interleaving of the sequence of steps performed by each process. All objects are *linearizable* or *atomic* in the sense that the processes obtain results from their operations on an object as if those operations were performed sequentially in the order specified by the execution. The *configuration* at any point in an execution is given by the state of all processes and the value of all objects. A process may become faulty at a given point in an execution, in which case it performs no subsequent operations.

A more formal description of our model can be given in terms of I/O automata [Lynch 1996; Lynch and Tuttle 1987; 1988; Pogosyants et al. 1996]. The randomized asynchronous computation model, which includes coin flips, is described in Aspnes and Herlihy [1990]. A formal definition of linearizability can be found in Herlihy and Wing [1990].

An operation of an object type is said to be *trivial* if applying the operation to any object of the type always leaves the value of the object unchanged. The READ operation is an example of a trivial operation. Two operations on an object type are said to *commute* if the order in which those operations are applied to any object of the type does not affect the resulting value of the object. For example, DECREMENT and FETCH&ADD operations commute with themselves and one another. A trivial operation commutes with any other operation on the same object.

An operation f on an object *overwrites* an operation f' if, starting from any value, performing f' and then f results in the same value (or set of values) as performing just f (i.e., $f(x) = f(f'(x))$ for all possible values x). This implies that, from every configuration, every sequence of operations on that object yields the same sequence of responses when preceded by f as when preceded by f' and f . If the value transition relation associated with an operation is idempotent, then the operation overwrites itself. All WRITE, TEST&SET, and SWAP operations on an object overwrite one another.

An object type is *historyless* if all its nontrivial operations overwrite one another. In other words, the value of a historyless object depends only on the last nontrivial operation that was applied to it.

A set of operations on an object is *interfering* if every pair of these operations either commute or overwrite one another. For example, the set of READ, WRITE, and SWAP operations is interfering, but the set of COMPARE&SWAP operations is not.

An *implementation* of an object X is a set of objects Y_1, \dots, Y_m representing X together with procedures F_1, \dots, F_n called by processes P_1, \dots, P_n to

execute operations on X . For object types X and Y , we say that X can be implemented from m instances of Y if there exists an implementation of an object of type X by processes P_1, \dots, P_n using objects Y_1, \dots, Y_m of type Y .

An implementation is *wait-free* if *each* nonfaulty process P_i always finishes executing its procedure F_i within a fixed, finite number of its own steps, regardless of the pace or failure of other processes. An implementation is *nonblocking* if, for every configuration and every execution starting at that configuration, there is *some* process P_i that finishes executing its procedure F_i within a finite number of steps. Wait-free implies nonblocking. The nonblocking property permits individual processes to starve, but requires the system as a whole to make progress. The wait-free property excludes starvation: any process that continues to execute events will finish its current operation. If each procedure F_i can be called only a finite number of times, then wait-free is the same as nonblocking.

The wait-free and nonblocking properties can be extended to *randomized wait-free* and *randomized nonblocking* by only requiring that P_i finishes executing its procedure F_i within a finite *expected* number of its steps. (See Aspnes and Herlihy [1990] for a more formal definition.) If each procedure F_i can be called only a finite number of times, then randomized wait-free is the same as randomized nonblocking. Furthermore, if an algorithm is wait-free, then it is randomized wait-free and if it is nonblocking, then it is randomized nonblocking.

A *solo execution* is an execution all of whose steps are performed by one process. An implementation has the *nondeterministic solo termination property* if, for every configuration C and every process P_i , there exists a finite solo execution, starting at configuration C , in which P_i finishes executing its procedure F_i . In other words, if P_i has no interference, it will finish performing its operation. If an algorithm is randomized wait-free (or wait-free), then it has the nondeterministic solo termination property, since every state transition having nonzero probability can be viewed as a possible nondeterministic choice.

Nondeterministic solo termination is a strictly weaker property than wait-freedom and randomized wait-freedom. For example, the simple snapshot algorithm following Observation 1 in Afek et al. [1993] is not randomized wait-free, but satisfies the nondeterministic solo termination property.

We evaluate the randomized space complexity of an object type based on the number of objects of the type that are required to implement n -process *binary consensus* in a randomized wait-free manner. A *binary consensus object* is an object on which each of n -processes can perform one DECIDE operation with input value in $\{0, 1\}$, returning an output value x , also in $\{0, 1\}$. The object's legal executions are those that satisfy the following conditions:

Consistency: The DECIDE operations of all processes return the same value.

Validity: If x is the value returned for some DECIDE operation, then x is the input value for the DECIDE operation of some process.

The first condition guarantees that consensus is achieved and the second condition excludes trivial solutions in which the outcome is fixed ahead of time.

We will use the term, *consensus*, to mean *n-process binary consensus*. A set of objects is said to solve *randomized consensus* if there is a randomized wait-free implementation of consensus from only that set of objects. No executions of an

implementation may give an incorrect answer (i.e., one that violates consistency or validity). In other words, we do not consider Monte Carlo implementations.

An execution of an implementation of consensus is *terminating* if it completes a DECIDE operation. Arbitrarily long and even nonterminating executions are possible in a randomized wait-free implementation. For example, since it is impossible to implement consensus in a wait-free manner for two or more processes from only read-write registers, any randomized wait-free implementation of consensus for two or more processes from only read-write registers must have nonterminating executions. However, these executions must occur with correspondingly small probabilities.

Consider the situation where there is a lower bound on the number of objects to solve randomized consensus, for objects of a particular type. Then, this bound can be used to obtain lower bounds on the number of instances of that object type that are necessary to implement objects of other types.

THEOREM 2.1. *Let X and Y be object types. Suppose $f(n)$ instances of X solve n -process randomized consensus and $g(n)$ instances of Y are required to solve n -process randomized consensus. Then, any randomized nonblocking implementation (and, hence, any randomized wait-free implementation) of X by Y for n processes requires $g(n)/f(n)$ instances of Y .*

PROOF. Suppose there exists a randomized nonblocking implementation of X using $h(n)$ instances of Y . Let \mathcal{A} denote a randomized wait-free implementation of consensus using $f(n)$ instances of X . Construct a new randomized wait-free implementation of consensus by replacing each instance of X in \mathcal{A} with a randomized nonblocking implementation using $h(n)$ instances of Y . In total, this implementation uses $f(n)h(n)$ instances of Y . Therefore, $f(n)h(n) \geq g(n)$ so $h(n) \geq g(n)/f(n)$. \square

3. Lower Bounds

In this section, we prove an $\Omega(\sqrt{n})$ lower bound on the number of objects required to solve randomized consensus, if each object is historyless. More specifically, we do this by showing that there is no implementation of consensus satisfying nondeterministic solo termination from a small number of historyless objects. The nondeterministic solo termination property is weaker than randomized wait-freedom, but it is sufficient for proving our lower bound.

An implementation of consensus is required to be consistent in every execution. Therefore, to demonstrate that an implementation of consensus is faulty, it suffices to exhibit an execution in which one process decides the value 0 and another process decides the value 1. This is done by combining an execution that decides 0 with an execution that decides 1. Although our lower bound is stated in terms of a strong model that requires objects to be linearizable, it also applies to weaker models, such as those which guarantee only sequential consistency.

Throughout this section, we use the following notation. The number of processes used in the implementation under consideration is n and the number of objects used is r . If V is a subset of these objects, then \bar{V} denotes the subset of these objects not in V . The sizes of V and \bar{V} are denoted by v and \bar{v} , respectively. Similarly, if \mathcal{P} is a subset of the n processes, then $\bar{\mathcal{P}}$ denotes the subset of processes not in \mathcal{P} .

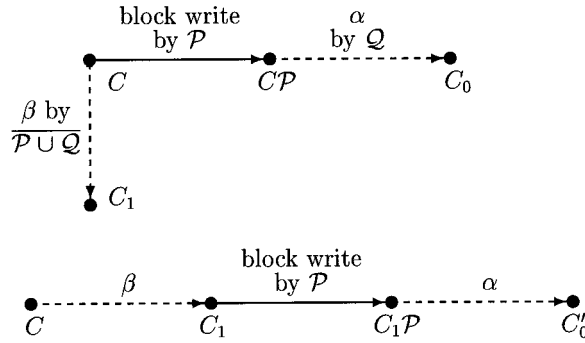


FIG. 1. Combining two executions.

In any configuration of the implementation, a process P is said to be *poised at* object R if P will perform a nontrivial (historyless) operation on R when next allocated a step. In this case, the value of object R can be fixed at any single future point in time by scheduling P to perform that operation. A *block write* to a set of objects V by a set of v processes \mathcal{P} consists of a sequence of v consecutive nontrivial (historyless) operations, one performed by each of the processes in \mathcal{P} to a different object in V . In any configuration C immediately before a block write to a set of objects V by a set of processes \mathcal{P} , there must be one process in \mathcal{P} poised at each object in V . The configuration that results from performing this block write will be denoted $C\mathcal{P}$. Note that a block write to V fixes the values of all the objects in V .

Under certain favorable conditions, it is possible to combine an execution that decides 0 with an execution that decides 1. For example, suppose there is a configuration C with the following properties: at C , there is a set \mathcal{P} of r processes that can perform a block write to the entire set of r objects; from $C\mathcal{P}$, there is an execution α by a disjoint set of processes \mathcal{Q} that leads to a configuration C_0 at which 0 is decided; and from C , there is an execution β containing steps of processes in neither \mathcal{P} nor \mathcal{Q} that leads to a configuration C_1 at which 1 is decided. Then the following is an execution from C that decides both 0 and 1. First perform β . Then let the processes in \mathcal{P} perform a block write. Finally, perform α . This is illustrated in Figure 1.

Note that the configurations $C\mathcal{P}$ and $C_1\mathcal{P}$ are indistinguishable to the processes in \mathcal{Q} . By writing, the processes in \mathcal{P} have obliterated all traces of β from the objects, so β is invisible to the processes in \mathcal{Q} .

3.1. READ WRITE REGISTERS AND IDENTICAL PROCESSES. First, we prove a lower bound in a much simpler situation: the objects are *read-write* registers (i.e., the only operations are READ and WRITE) and all processes are identical. Thus, if two processes are in the same state, they perform the same operation on the same register when they are next allocated a step and, if the outcomes of those operations are the same (e.g., the values of their coin flips or the values that they read), their resulting states will be the same. Furthermore, processes with the same input value will be in the same initial state. Although the lower bound proof in this restricted setting is considerably easier than in the general case, the overall structure of both proofs are similar and we feel the easier proof provides important intuition.

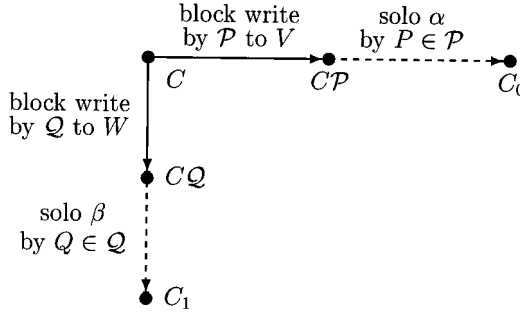


FIG. 2. Conditions of Lemma 3.1.

One technique that can only be applied in the restricted setting is *cloning*. Consider any point in some execution at which a process P writes a value x to a register R . Then there is another execution which is the same except that a group of clones have been left behind, all poised to write value x to register R . To construct this execution, the clones are given the same initial state as P and then P and its clones are scheduled as a group, up to the point at which P performs the designated write of x to R . In other words, whenever P is allocated a step, each of its clones is immediately allocated a step. The outcomes of the clones' internal operations in this execution are specified to be the same as the corresponding outcomes of P 's internal operations. Then, up to the point at which P performs the write, the clones have the same state as P and they perform exactly the same sequence of operations. These clones can be scheduled to perform their writes at various subsequent points of time, resetting the contents of the register R to the same value x each time.

Provided there are sufficiently many processes available, cloning enables two executions to be combined into an inconsistent execution under the general conditions illustrated in Figure 2.

LEMMA 3.1. Consider any implementation of consensus from r read-write registers using identical processes that satisfies nondeterministic solo termination. Let C be a configuration in which there is a set of $v \geq 1$ processes \mathcal{P} poised at some set of registers V and a disjoint set of $w \geq 1$ processes \mathcal{Q} poised at some (not necessarily disjoint) set of registers W . Suppose that, after the block write to V by the processes in \mathcal{P} , there is a solo execution α by a process P in \mathcal{P} that decides 0 and, symmetrically, after the block write to W by the processes in \mathcal{Q} , there is a solo execution β by a process Q in \mathcal{Q} that decides 1. Then there is an execution from C that decides both 0 and 1 and uses at most $r^2 - r + (3v + 3w - v^2 - w^2)/2$ identical processes.

PROOF. The proof is by induction on $\bar{v} + \bar{w}$.

Case 1. $V \subseteq W$. (Note that this is the case when $\bar{w} = 0$.)

If all writes in α are to registers in W , consider the following execution starting from C . It is illustrated in Figure 3. First, processes in \mathcal{P} perform a block write to V , next α is performed (taking the system to configuration C_0), and then processes in \mathcal{Q} perform a block write to W . Note that the resulting configuration $C_0\mathcal{Q}$ is indistinguishable to processes in \mathcal{Q} from the configuration $C\mathcal{Q}$ obtained from C by performing the block write to W . Finally, β is performed. This

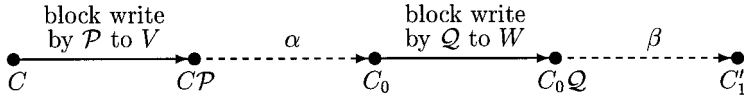


FIG. 3. $V \subseteq W$ and α writes only to W .

execution decides both 0 and 1 and uses $v + w$ processes. Since $v, w \leq r$, it follows that $v + w \leq r^2 - r + (3v + 3w - v^2 - w^2)/2$.

Otherwise, there is at least one a write to a register not in W during the solo execution α . Let C' be the configuration just before the first such write occurs, let $R \notin W$ be the register written to by that write, let α' denote that portion of α occurring after the write, and let $V' = V \cup \{R\}$. Note that $R \notin V$, because $R \notin W$ and $V \subseteq W$. Hence, $v' = v + 1$. During the execution from C to C' , every register in V is written to at least once. Thus, if there are sufficiently many processes, a clone can be left poised at each register in V , ready to re-perform the last write that was performed on the register prior to C' . These v clones, together with the process P performing the solo execution α (which includes the write to R and the solo execution α'), will form \mathcal{P}' . Note that the value of every register is the same in the configurations $C' \setminus \{P\}$ and $C' \mathcal{P}'$. Since α' decides the value 0 starting at $C' \setminus \{P\}$, it is also the case that α' decides 0 starting at $C' \mathcal{P}'$. Furthermore, since $V \subseteq W$, the configurations C and C' are indistinguishable to processes in \mathcal{Q} . Therefore, β decides 1 starting at $C' \mathcal{Q}$. This is illustrated in Figure 4. By the induction hypothesis, there is an execution from C' that decides both 0 and 1 and uses at most $r^2 - r + (3(v + 1) + 3w - (v + 1)^2 - w^2)/2$ processes. Prepending the execution from C to C' yields an execution from C that decides both 0 and 1 and uses $(v - 1)$ additional processes (those in $\mathcal{P} - \mathcal{P}'$) for a total of at most $r^2 - r + (3v + 3w - v^2 - w^2)/2$.

Case 2. $W \subseteq V$. By symmetry, there is an execution starting from C that decides both 0 and 1.

Case 3. $V \not\subseteq W$ and $W \not\subseteq V$. Consider any terminating execution from C that begins with a block write to $U = V \cup W$ and continues with a solo execution γ by one of these u processes. Such an execution exists by the nondeterministic solo termination property. Since $U \subseteq V$ and $V \neq \emptyset$, it follows that $u - 1 \geq v \geq 1$. Suppose γ decides 0. (The case when γ decides 1 is symmetric.) Let \mathcal{P}'' consist of \mathcal{P} plus a clone of each process in \mathcal{Q} that is poised at a register in $W - V$ in configuration C . This is illustrated in Figure 5.

By the induction hypothesis applied to configuration C , sets of registers U and W , and the sets of processes \mathcal{P}' and \mathcal{Q} , there is an execution from C which decides both 0 and 1 and uses at most $r^2 - r + (3u + 3w - u^2 - w^2)/2 \leq r^2 - r + (3v + 3w - v^2 - w^2)/2$ processes. \square

To obtain our lower bound, we show that, for any randomized consensus algorithm, the conditions necessary to apply Lemma 3.1 can be achieved, provided sufficiently many processes are available.

LEMMA 3.2. *There is no implementation of consensus satisfying nondeterministic solo termination from r read-write registers using $r^2 - r + 2$ or more identical processes.*

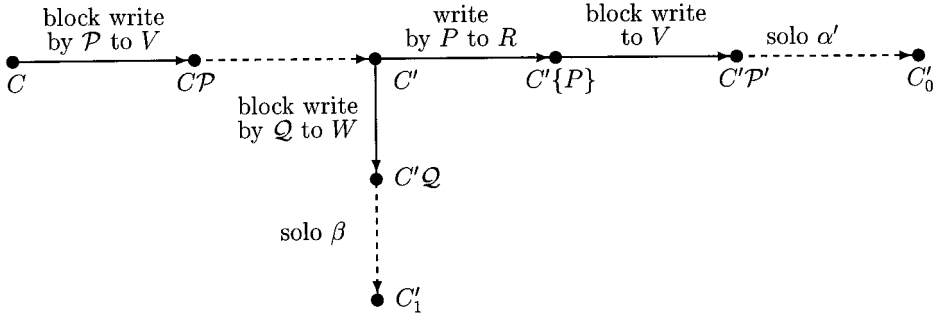


FIG. 4. $V \subseteq W$ and α writes to $R \notin W$.

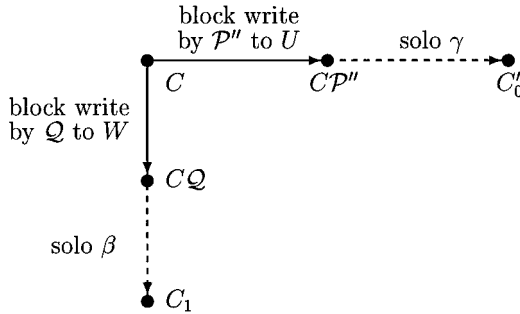


FIG. 5. $V \notin W, W \notin V$, and γ decides 0.

PROOF. Consider any implementation of consensus satisfying nondeterministic solo termination from r read-write registers using $r^2 - r + 2$ or more identical processes.

Let P and Q be processes with initial values 0 and 1, respectively. Let α denote any terminating solo execution by P and let β denote any terminating solo execution by Q . The existence of these executions is guaranteed by the nondeterministic solo termination property. Furthermore, by validity, α must decide 0 and β must decide 1.

If one of these executions, say α , contains no writes, then the execution consisting of α followed by β decides both 0 and 1. Therefore, we may assume that both α and β contain at least one write.

Let α' denote the portion of α occurring after the first write and let V be the singleton set consisting of the register P first writes to. Define β' and W analogously. Let γ be the execution consisting of those operations of α and β occurring before their first writes and let C be the configuration obtained from the initial configuration by performing γ . If there are at least $r^2 - r + 2$ processes, it follows by Lemma 3.1 that there is an execution from C that decides both 0 and 1. Prepending this execution by γ yields an execution from the initial configuration that decides both 0 and 1. This violates the consistency condition. \square

The following result is a direct consequence of Lemma 3.2.

THEOREM 3.3. *A randomized wait-free implementation of binary consensus for n identical processes using only read-write registers requires $\Omega(n)$ registers.*

3.2. GENERAL CASE. Next, we show our main result: $\Omega(\sqrt{n})$ objects are necessary to implement n -process binary consensus in a randomized wait-free manner, if the objects are historyless. The key to the lower bound is the definition of an interruptible execution. Informally, an interruptible execution is an execution that can be broken into pieces, between which executions involving other processes can be inserted.

Each piece of an interruptible execution begins with a block write to a set of objects. Because the objects are historyless, this block write fixes these objects to have particular values, no matter when it is executed. Hence, if an execution performed by other processes only changes the values of objects in the set, then that execution can be inserted immediately before this piece without affecting the rest of the interruptible execution.

Note that, unlike the special case considered in Section 3.1, the resulting state of each process performing the block write may depend on the value of the object it accesses. This could influence a subsequent execution involving this process. Therefore, in the definition of an interruptible execution, processes that participate in a block write take no further steps.

Definition 3.1. An execution α starting from configuration C is *interruptible* with *initial object set* V and *process set* \mathcal{P} if all steps in α are taken by processes in \mathcal{P} and α can be divided into one or more *pieces* $\alpha = \alpha_1 \cdots \alpha_k$ such that

- α_i begins with a block write to a set of objects V_i by processes that take no further steps in α ,
- all nontrivial operations in α_i are to objects in V_i ,
- $V = V_1 \not\subseteq \cdots \not\subseteq V_k$, and
- after α has been performed, some process has decided.

In other words, if an execution α is interruptible with initial object set V then $\alpha = \alpha_1\alpha'$, where α_1 begins with a block write to V by processes that take no further steps in α and α_1 contains no nontrivial operations on objects in \bar{V} . Furthermore, if C' is the configuration obtained from C by executing α_1 , then either some process has decided at C' and α' is empty or α' is an interruptible execution starting from C' with some initial object set $V' \not\subseteq V$.

We say that an interruptible execution *decides* a value x if, after α has been performed, some process has decided x .

A terminating solo execution augmented with a sufficient numbers of clones performing block writes at appropriate points is a special case of an interruptible execution. This is what is used to obtain the lower bound for identical processes. However, when processes are not assumed to be identical, clones cannot necessarily be created. The processes performing the block writes at the beginning of the pieces of an interruptible execution are used instead of clones in many places throughout the following lower bound proof. The excess capacity of an interruptible execution, defined below, is used to satisfy an additional need for clones.

Definition 3.2. An interruptible execution $\alpha = \alpha_1 \cdots \alpha_k$ starting from configuration C with initial object set V and process set \mathcal{P} has *excess capacity e for object set U* if, at the beginning of each piece α_i , there are at least e processes not in \mathcal{P} poised at each object in $V_i \cap U$.

Note that if $k > 1$ and C' is the configuration obtained from C by executing α_1 , then interruptible execution $\alpha_2 \cdots \alpha_k$ also has excess capacity e for U .

The next lemma shows that, from any configuration in which there are sufficiently many processes poised at certain objects, there is an interruptible execution with desired excess capacity. Note that there was no need to prove an analogous lemma in Section 3.1, since the nondeterministic solo termination property guarantees the existence of terminating solo executions (the analogue of interruptible executions).

LEMMA 3.4. *Let U and V be sets of objects and let \mathcal{P} be a set of at least $(r^2 + r - v^2 + v)/2 + e|\bar{V} \cap U|$ processes. Consider any configuration C in which there are at least $\bar{v} + 1$ processes in \mathcal{P} poised at every object in V and at least e processes not in \mathcal{P} poised at every object in $V \cap U$. Then there is an interruptible execution starting from C with initial object set V and process set \mathcal{P} that has excess capacity e for U .*

PROOF. By induction on \bar{v} . Consider any execution δ starting from configuration C with the following properties:

- there is a set $\hat{\mathcal{P}} \subseteq \mathcal{P}$ of $v(\bar{v} + 1)$ processes which, at C , contains $\bar{v} + 1$ processes poised at each object in V ,
- δ begins with a block write to V by a set $\mathcal{P}_1 \subseteq \hat{\mathcal{P}}$ of v processes that take no further steps in δ ,
- all other steps in δ are taken by processes in $\mathcal{P} - \hat{\mathcal{P}}$,
- all nontrivial operations in δ are performed on objects in V , and
- at the configuration C' obtained from C by executing δ , either some process has decided or all processes in $\mathcal{P} - \hat{\mathcal{P}}$ are poised at objects in \bar{V} .

Such an execution may be obtained by performing the block write to V by \mathcal{P}_1 and then, for each process in $\mathcal{P} - \hat{\mathcal{P}}$, performing steps of a solo terminating execution until the process has decided or is poised at an object in \bar{V} .

If, at C' , some process has decided (which must be the case if $\bar{v} = 0$), then δ is an interruptible execution with one piece that satisfies the desired conditions. Therefore, assume that $\bar{v} > 0$ and, at C' , every process in $\mathcal{P} - \hat{\mathcal{P}}$ is poised at an object in \bar{V} .

For every integer $i \geq 1$, let y_i and z_i be the number of objects in $\bar{V} \cap \bar{U}$ and $\bar{V} \cap U$, respectively, at which at least i processes in $\mathcal{P} - \hat{\mathcal{P}}$ are poised in configuration C' . Then $y_i \geq y_{i+1}$ and $z_i \geq z_{i+1}$.

Furthermore, there exists $i \in \{1, \dots, \bar{v}\}$ such that $y_i + z_{e+i} \geq \bar{v} - i + 1$. To see why, suppose to the contrary that $y_i + z_{e+i} \leq \bar{v} - i$ for all $1 \leq i \leq \bar{v}$. In particular, $y_{\bar{v}} + z_{e+\bar{v}} \leq 0$, so $y_i = z_{e+i} = 0$ for all $i \geq \bar{v}$. Then, since $v + \bar{v} = r$ and $\bar{v} > 0$,

$$\begin{aligned}
 |\mathcal{P} - \hat{\mathcal{P}}| &= \sum_{i \geq 1} (y_i + z_i) \\
 &= \sum_{i=1}^{\bar{v}-1} (y_i + z_{e+i}) + \sum_{i=1}^e z_i
 \end{aligned}$$

$$\begin{aligned}
 &\leq \sum_{i=1}^{\bar{v}-1} (\bar{v} - i) + \sum_{i=1}^e |\bar{V} \cap U| \\
 &= \frac{\bar{v}(\bar{v} - 1)}{2} + e|\bar{V} \cap U| \\
 &\leq \frac{\bar{v}(\bar{v} - 1)}{2} + |\mathcal{P}| - \frac{(r^2 + r - v^2 + v)}{2} \\
 &= |\mathcal{P}| - v\bar{v} - v - \bar{v} \\
 &< |\mathcal{P}| - v(\bar{v} + 1) \\
 &= |\mathcal{P} - \hat{\mathcal{P}}|.
 \end{aligned}$$

This is a contradiction.

Consider the situation at configuration C' . Suppose $Y \subseteq \bar{V} \cap \bar{U}$ and $Z \subseteq \bar{V} \cap U$ are sets of objects such that there are at least i processes in $\mathcal{P} - \hat{\mathcal{P}}$ poised at every object in Y , there are at least $e + i$ processes in $\mathcal{P} - \hat{\mathcal{P}}$ poised at every object in Z , and $|Y| + |Z| = \bar{v} - i + 1$. Let $V' = V \cup Y \cup Z$. Then $v' = v + \bar{v} - i + 1 = r - i + 1$, so $i = \bar{v}' + 1$.

Let $\mathcal{E} \subseteq \mathcal{P} - \hat{\mathcal{P}}$ be a set of $e|Z|$ processes, e poised at every object in $Z = (V' - V) \cap U$. Define $\mathcal{P}' = \mathcal{P} - \mathcal{P}_1 - \mathcal{E}$. At C , there are at least e processes not in \mathcal{P} (and hence not in \mathcal{P}') poised at every object in $V \cap U$ and these processes take no steps in δ . Also, none of the processes in \mathcal{E} are in \mathcal{P}' . Therefore, at C' , there are at least e processes not in \mathcal{P}' poised at every object in $V' \cap U$.

There are at least $i = \bar{v}' + 1$ processes in $\mathcal{P} - \hat{\mathcal{P}} - \mathcal{E} \subseteq \mathcal{P}'$ poised at every object in $Y \cup Z = V' - V$. Furthermore, since none of the processes in $\hat{\mathcal{P}} - \mathcal{P}_1$ take any steps in δ , it follows that, at configuration C' , there are also at least $\bar{v} \geq \bar{v}' + 1$ processes in \mathcal{P}' poised at every object in V .

Finally, since $v \leq v' - 1$,

$$\begin{aligned}
 |\mathcal{P}'| &= |\mathcal{P}| - v - e|(V' - V) \cap U| \\
 &\geq \frac{r^2 + r - v^2 + v}{2} + e|\bar{V} \cap U| - v - e|(V' - V) \cap U| \\
 &= \frac{r^2 + r - v^2 - v}{2} + e|\bar{V}' \cap U| \\
 &\geq \frac{r^2 + r - (v')^2 + v'}{2} + e|\bar{V}' \cap U|.
 \end{aligned}$$

Then, by the induction hypothesis, there is an interruptible execution δ' starting from C' with initial object set V' and process set \mathcal{P}' that has excess capacity e for U .

Now $\mathcal{P}' \subseteq \mathcal{P} - \mathcal{P}_1$ and $V' \not\subseteq V$. Therefore, $\delta\delta'$ is an interruptible execution starting from C with initial object set V and process set \mathcal{P} that has excess capacity e for U . \square

The next result describes conditions under which two interruptible executions can be combined to form an inconsistent execution. It is analogous to Lemma 3.1, but is more difficult, because we cannot just make a sufficient number of clones as we need them. Instead, we use excess capacity in one execution to guarantee a sufficient number of processes poised at the particular objects needed for the other execution.

LEMMA 3.5. *Let α be an interruptible execution, starting at configuration C , with initial object set V , process set \mathcal{P} , and excess capacity \bar{w} for \bar{W} and let β be an interruptible execution, starting at configuration C , with initial object set W , process set \mathcal{Q} , and excess capacity \bar{v} for \bar{V} . Suppose $|\mathcal{P}| \geq (r^2 + r - v^2 + v)/2 + \bar{w} \cdot |\bar{V} \cap \bar{W}|$, $|\mathcal{Q}| \geq (r^2 + r - w^2 + w)/2 + \bar{v} \cdot |\bar{V} \cap \bar{W}|$, and \mathcal{P} and \mathcal{Q} are disjoint. If α decides 0 and β decides 1, then there is an execution starting from C that decides both 0 and 1.*

PROOF. By induction on $\bar{v} + \bar{w}$.

Case 1. $V \subseteq W$. (Note that this is the case when $\bar{w} = 0$).

Let α_1 be the first piece of α and let C' be the configuration obtained from C by executing α_1 . All of the nontrivial operations in α_1 are to objects in $V \subseteq W$, \mathcal{P} and \mathcal{Q} are disjoint, and β begins with a block write to W by a set \mathcal{Q}_1 of w processes that take no further steps in β . This implies that the configurations $C\mathcal{Q}_1$ and $C'\mathcal{Q}_1$ are indistinguishable to the processes in $\mathcal{Q} - \mathcal{Q}_1$. Furthermore, all processes poised at objects in \bar{V} at configuration C are poised at the same objects at configuration C' . Therefore, β is an interruptible execution starting from C' , with initial object set W and process set \mathcal{Q} , that decides 1 and has excess capacity \bar{v} for \bar{V} . If some process has decided 0 at C' , then $\alpha_1\beta$ is an execution from C that decides both 0 and 1. Otherwise, $\alpha = \alpha_1\alpha'$, where α' is an interruptible execution α' starting from C' with process set \mathcal{P} and some initial object set $V' \not\subseteq V$ that decides 0 and has excess capacity \bar{w} for \bar{W} . Since $\bar{V}' \not\subseteq \bar{V}$,

$$|\mathcal{P}| \geq \frac{r^2 + r - v^2 + v}{2} + \bar{w}|\bar{V} \cap \bar{W}| \geq \frac{r^2 + r - (v')^2 + v'}{2} + \bar{w}|\bar{V}' \cap \bar{W}|,$$

$$|\mathcal{Q}| \geq \frac{r^2 + r - w^2 + w}{2} + \bar{v}|\bar{V} \cap \bar{W}| > \frac{r^2 + r - w^2 + w}{2} + \bar{v}'|\bar{V}' \cap \bar{W}|,$$

and β has excess capacity \bar{v}' for \bar{V}' . By the induction hypothesis applied to α' and β , there is an execution δ starting from C' that decides both 0 and 1. Hence, $\alpha_1\delta$ is an execution starting from C that decides both 0 and 1.

Case 2. $W \subseteq V$. Similarly, there is an execution starting from C that decides both 0 and 1.

Case 3. $V \not\subseteq W$ and $W \not\subseteq V$. Let $V' = W' = V \cup W$. Consider the situation at configuration C . Since β has excess capacity $\bar{v} \geq \bar{v}' + 1$ for $\bar{V} \supseteq \bar{V}'$, β has excess capacity \bar{v}' for \bar{V}' . In addition, there are $\bar{v}' + 1$ processes not in \mathcal{Q} poised

at each object in $W \cap \bar{V} = V' - V$. Form \mathcal{P}' by adding these processes to \mathcal{P} , if they are not already in \mathcal{P} . Note that \mathcal{P}' and \mathcal{Q} are disjoint. There are also at least $\bar{v} + 1 \geq \bar{v}' + 1$ processes in $\mathcal{P} \subseteq \mathcal{P}'$ poised at each object in V , by assumption. Since α has excess capacity \bar{w} for \bar{W} , there are \bar{w} processes not in \mathcal{P} poised at each object in $V \cap \bar{W} = V' \cap \bar{W}$. These processes are not in \mathcal{P}' , because the processes in $\mathcal{P}' - \mathcal{P}$ are poised at objects in \bar{V} . Since $V \not\subseteq V'$,

$$\begin{aligned} |\mathcal{P}'| &\geq |\mathcal{P}| \geq \frac{r^2 + r - v^2 + v}{2} + \bar{w} \cdot |\bar{V} \cap \bar{W}| \\ &\geq \frac{r^2 + r - (v')^2 + v'}{2} + \bar{w} \cdot |\bar{V}' \cap \bar{W}|. \end{aligned}$$

By Lemma 3.4, there exists an interruptible execution α' starting from configuration C with initial object set V' and process set \mathcal{P}' that has excess capacity \bar{w} for \bar{W} . If α' decides 0, then it follows from the induction hypothesis applied to α' and β that there is an execution starting from C that decides both 0 and 1. Therefore we may assume that α' decides 1. Note that α' has excess capacity \bar{w}' for \bar{W}' , since $\bar{W}' \not\subseteq \bar{W}$.

Similarly, we may assume that there exists an interruptible execution β' starting from C with initial object set W' and process set \mathcal{Q}' that decides 0 and has excess capacity \bar{v} for \bar{V} , where $|\mathcal{Q}'| \geq (r^2 + r - (w')^2 + w')/2 + \bar{v} \cdot |\bar{W}' \cap \bar{V}|$, $\mathcal{Q}' \supseteq \mathcal{Q}$ is disjoint from \mathcal{P} , and all processes in $\mathcal{Q}' - \mathcal{Q}$ are poised at objects in $W' - W$.

Since $(V' - V)$ and $(W' - W)$ are disjoint, $\mathcal{P}' - \mathcal{P}$ and $\mathcal{Q}' - \mathcal{Q}$ are disjoint. By assumption, \mathcal{P} and \mathcal{Q} are disjoint. It follows that \mathcal{P}' and \mathcal{Q}' are disjoint. Furthermore, $V' \not\supseteq V$ and $W' \not\supseteq W$, so $|\mathcal{P}'| > (r^2 - r - (v')^2 + v')/2 + \bar{w}' \cdot |\bar{V}' \cap \bar{W}'|$ and $|\mathcal{Q}'| > (r^2 - r - (w')^2 + w')/2 + \bar{v}' \cdot |\bar{V}' \cap \bar{W}'|$. Then, by the induction hypothesis applied to β' and α' , there is an execution starting from C that decides both 0 and 1. \square

LEMMA 3.6. *There is no implementation of consensus satisfying nondeterministic solo termination from r historyless objects using $3r^2 + r$ or more processes.*

PROOF. Consider any (randomized) algorithm that purports to achieve wait-free binary consensus among $3r^2 + r$ processes using r objects. Partition these processes into two sets, \mathcal{P} and \mathcal{Q} , each containing $(3r^2 + r)/2$ processes. Give each process in \mathcal{P} the initial value 0 and give each process in \mathcal{Q} the initial value 1.

Let $V = W = \emptyset$. By Lemma 3.4, there is an interruptible execution α starting from the initial configuration with initial object set V and process set \mathcal{P} that has excess capacity \bar{w} for \bar{W} . Since the processes in \mathcal{P} all have initial value 0, α must decide 0. Similarly, there is an interruptible execution β starting from the initial configuration with initial object set W and process set \mathcal{Q} that has excess capacity \bar{v} for \bar{V} and decides 1. Hence, Lemma 3.5 implies that there is an execution starting from the initial configuration that decides both 0 and 1, violating the consistency condition. \square

The following result is a direct consequence of Lemma 3.6.

THEOREM 3.7. *A randomized wait-free implementation of n -process binary consensus requires $\Omega(\sqrt{n})$ objects, if the objects are historyless.*

4. Separation Results

We use our main theorem to derive a series of results comparing the “randomized power” of various synchronization primitives with their “deterministic power” [Herlihy 1991]. We say that object type X is *deterministically more powerful* than object type Y if the number of processes for which consensus can be achieved is larger using instances of X than using instances of Y . For randomized computation, we say that object type X is *more powerful* than object type Y if n -process randomized consensus requires asymptotically fewer instances of X than instances of Y .

Objects with only interfering operations cannot deterministically solve 3-process consensus [Herlihy 1991]. Hence, they are deterministically less powerful than objects with operations such as COMPARE&SWAP, which can solve n -process consensus.

Consider any object with an operation such that, starting with some particular state, the response from one application of the operation is always different than the response from the second of two successive applications of that operation. (For example, a register with the value 0 returns different values from successive applications of SWAP(1). The operation FETCH&ADD applied starting with any value also has this property.) Then this object can solve 2-process consensus. Therefore, it is deterministically more powerful than the *read-write* register, which cannot solve 2-process consensus.

Herlihy [1991, Theorem 5] shows that n -process consensus can be implemented deterministically using a single bounded *compare&swap* register. Hence, from Theorems 2.1 and 3.7, we have the following result.

COROLLARY 4.1. *Any randomized nonblocking bounded compare&swap register implementation requires $\Omega(\sqrt{n})$ objects, if the objects are historyless.*

Since there are deterministic counter implementations using $O(n)$ read-write registers [Aspnes and Herlihy 1990; Moran et al. 1992], neither counters nor bounded counters can deterministically solve 2-process consensus [Herlihy 1991].

Aspnes [1990] gives a randomized algorithm for n process binary consensus using three bounded counters: the first two keep track of the number of processes with input 0 and input 1 respectively, and the third is used as the cursor for a random walk. The first two counters assume values between 0 and n , while the third assumes values between $-3n$ and $3n$. The first two counters can be eliminated at some cost in performance (J. Aspnes, private communication).

THEOREM 4.2 (ASPNES). *There is a randomized consensus implementation using one bounded counter.*

The next result follows from Theorems 2.1, 3.7, and 4.2.

COROLLARY 4.3. *Any randomized nonblocking bounded counter implementation requires $\Omega(\sqrt{n})$ objects, if the objects are historyless.*

Surprisingly, the lower bounds in Corollaries 4.1 and 4.3 are both independent of the number of values an object can assume: they hold even when the objects (such as *read-write* and *swap* registers) used in the implementation have an infinite number of values, but the object being implemented has a finite number of values.

The same result holds for the implementation of a *fetch&add* register, a *fetch&increment* register, or a *fetch&decrement* register, because a single instance of any of these objects can be easily used to implement a counter.

THEOREM 4.4. *Randomized consensus can be solved using a single instance of a fetch&add register, a fetch&increment register, or a fetch&decrement register.*

COROLLARY 4.5. *Any randomized nonblocking implementation of a fetch&add register, a fetch&increment register, or a fetch&decrement register requires $\Omega(\sqrt{n})$ objects, if the objects are historyless.*

Theorem 4.4 is particularly interesting since it shows that *fetch&add* and *compare&swap*, which differ substantially in their deterministic power [Herlihy 1991], have similar randomized power in the sense that one instance of each suffices to solve randomized consensus.

5. Conclusions

We have presented a separation among multiprocessor synchronization primitives based on the *space complexity* of randomized solutions to n -process binary consensus. Our main result proved that $\Omega(\sqrt{n})$ objects are necessary to solve randomized n -process binary consensus, if the objects are historyless. Randomized n -process consensus can be solved using $O(n)$ read-write registers of bounded size and we conjecture that the true space complexity of this problem is $\Theta(n)$. We believe that, based on our approach, a larger lower bound may be possible by reusing the processes that perform the block writes at the beginning of a piece of an interruptible execution.

We further believe that our lower bound approach can be extended to allow comparisons among other classes of primitives, and help us to better understand the limitations of using randomization to implement various synchronization primitives from one another. The lower bounds presented here only consider the implementation of a “single access” object, but also apply to the implementation of a “multiple use” object, where each process can access the object repeatedly. However, it may be that the implementation of certain multiple use objects, for example, real-world synchronization primitives such as *test&set* and *fetch&add*, is significantly more difficult and that improved lower bounds can be obtained by having some processors access the implemented object many times. Indeed, a recent result by Jayanti et al. [1996] shows that for multiple use objects, it takes $n - 1$ instances of objects such as registers or *swap* registers to implement objects such as *increment* registers, *fetch&add* registers, and *compare&swap* registers.

Needless to say, there is also much work to be done in providing efficient upper bounds for randomized implementations of objects.

REFERENCES

- ABRAHAMSON, K. 1988. On achieving consensus using a shared memory. In *Proceedings of the 7th Annual ACM Symposium on Principles of Distributed Computing* (Toronto, Ont., Canada, Aug. 15–17). ACM, New York, pp. 291–302.
- AFEK, Y., AND STUPP, G. 1993. Synchronization power depends on the register size. In *Proceedings of the 34th Annual IEEE Symposium on Foundations of Computer Science* (Nov.). IEEE Computer Science Press, Los Alamitos, Calif., pp. 196–205.

- AFEK, Y., ATTIYA, H., DOLEV, D., GAFNI, E., MERRITT, M., AND SHAVIT, N. 1993. Atomic snapshots of shared memory. *J. ACM* 40, 4 (Sept.), 873–890.
- AFEK, Y., GREENBERG, D., MERRITT, M., AND TAUBENFELD, G. 1995. Computing with faulty shared objects. *J. ACM* 42, 6 (Nov.), 1231–1274.
- ALEMANY, J., AND FELTEN, E. W. 1992. Performance issues in non-blocking synchronization on shared-memory multiprocessors. In *Proceedings of the 11th Annual ACM Symposium on Principles of Distributed Computing* (Vancouver, B.C., Canada, Aug. 10–12). ACM, New York, pp. 125–134.
- ASPINES, J. 1990. Time- and-space efficient randomized consensus. In *Proceedings of the 9th Annual ACM Symposium on Principles of Distributed Computing* (Quebec City, Que., Canada, Aug. 22–24). ACM, New York, pp. 325–331.
- ASPINES, J. 1997. Lower bounds for distributed coin-flipping and randomized consensus. In *Proceedings of the 29th Annual ACM Symposium on Theory of Computing* (El Paso, Tex., May 4–6). ACM, New York, pp. 559–568.
- ASPINES, J., AND HERLIHY, M. 1990. Fast, randomized consensus using shared memory. *J. Algorithms* 11 (Sept.), 441–461.
- ASPINES, J., AND WAARTS, O. 1992. Randomized Consensus in Expected $O(n \log^2 n)$ operations per processor. In *Proceedings of the 33rd Annual IEEE Symposium on the Foundation of Computer Science* (Oct.). IEEE Computer Science Press, Los Alamitos, Calif., pp. 137–146.
- ATTIYA, H., DOLEV, D., AND SHAVIT, N. 1989. Bounded polynomial randomized consensus. In *Proceedings of the 8th Annual ACM Symposium on Principles of Distributed Computing* (Edmonton, Alb., Canada, Aug. 14–16). ACM, New York, pp. 281–293.
- BERSHAD, B. 1993. Practical considerations for non-blocking concurrent objects. In *Proceedings of the 13th International Conference on Distributed Computing Systems* (May). IEEE Computer Society Press, Los Alamitos, Calif., pp. 264–274.
- BRACHA, G., AND RACHMAN, O. 1991. Randomized consensus in expected $O(n^2 \log n)$ operations. In *Proceedings of the 5th International Workshop on Distributed Algorithms* (Delphi, Greece, Oct.). Lecture Notes in Computer Science, vol. 579, Springer-Verlag, New York, pp. 143–150.
- BURNS, J., AND LYNCH, N. 1989. Mutual exclusion using indivisible reads and writes. In *Proceedings of the 18th Annual Allerton Conference on Communication Control, and Computing*. Monticello, Ill. pp. 833–842.
- CHOR, B., ISRAELI, A., AND LI, M. 1987. On processor coordination using asynchronous hardware. In *Proceedings of the 6th ACM Symposium on Principles of Distributed Computing* (Vancouver, B.C., Canada, Aug. 10–12). ACM, New York, pp. 86–97.
- DOLEV, D., DWORK, C., AND STOCKMEYER, L. 1987. On the minimal synchronism needed for distributed consensus. *J. ACM* 34, 1 (Jan.), 77–97.
- DWORK, C., HERLIHY, M., PLOTKIN, S. A., AND WAARTS, O. 1992. Time-lapse snapshots. Tech. Rep. STAN//CS-TR-92-1423. Dept. Computer Science, Stanford Univ., Stanford, Calif.
- GRAHAM, R. L., AND YAO, A. C. 1989. On the improbability of reaching Byzantine agreements. In *Proceedings of the 21st Annual ACM Symposium on the Theory of Computing* (Seattle, Wash., May 15–17). ACM, New York, pp. 467–478.
- HERLIHY, M. P. 1991a. Randomized wait-free concurrent objects. In *Proceedings of the 10th Annual ACM Symposium on Principles of Distributed Computing* (Montreal, Que., Canada, Aug. 19–21). ACM, New York, pp. 11–21.
- HERLIHY, M. P. 1991b. Wait-free synchronization. *ACM Trans. Prog. Lang. Syst.* 13, 1 (Jan.), 124–149.
- HERLIHY, M. P., AND WING, J. M. 1990. Linearizability: A correctness condition for concurrent objects. *ACM Trans. Prog. Lang. Syst.* 12, 3 (July), 463–492.
- JAYANTI, P., TAN, K., AND TOUEG, S. 1996. Time and space lower bounds for non-blocking implementations (preliminary version). In *Proceedings of the 15th Annual ACM Symposium on Principles of Distributed Computing* (Philadelphia, Pa., May 23–26). ACM, New York, pp. 257–266.
- KUSHILEVITZ, E., MANSOUR, Y., RABIN, M., AND ZUCKERMAN, D. 1993. Lower bounds for randomized mutual exclusion. In *Proceedings of the 25th Annual ACM Symposium on the Theory of Computing* (San Diego, Calif., May 16–18). ACM, New York, pp. 154–163.
- LAMPOR, L. 1986. On interprocess communication. Part II: Algorithms. *Dist. Comput.* 1, 2, 86–101.
- LOUI, M., AND ABU-AMARA, H. 1987. Memory requirements for agreement among unreliable asynchronous processes. *Adv. Comput. Res.* 4, 163–183.
- LYNCH, N. A. 1996. *Distributed Algorithms*. Morgan-Kaufmann, San Francisco, Calif.

- LYNCH, N. A., AND TUTTLE, M. R. 1987. Hierarchical correctness proofs for distributed algorithms. In *Proceedings of the 6th Annual ACM Symposium on Principles of Distributed Computing* (Vancouver, B.C., Canada, Aug. 10–12). ACM, New York, pp. 137–151. (Full version available as MIT Tech Rep. MIT/LCS/TR-387.)
- LYNCH, N. A., AND TUTTLE, M. R. 1988. An introduction to input/output automata. MIT Tech. Rep. MIT/LCS/TR-373.
- MORAN, S., TAUBENFELD, G., AND YADIN, I. 1992. Concurrent counting. In *Proceedings of the 11th Annual ACM Symposium on Principles of Distributed Computing* (Vancouver, B.C., Canada, Aug. 10–12). ACM, New York, pp. 59–70.
- POGOSYANTS, A., SEGALA, R., AND LYNCH, N. 1996. Verification of the randomized consensus algorithm of Aspnes and Herlihy: A case study. Unpublished manuscript. MIT, Cambridge, Mass.
- SAKS, M., SHAVIT, N., AND WOLL, H. 1991. Optimal time randomized consensus—Making resilient algorithms fast in practice. In *Proceedings of the 2nd Annual ACM Symposium on Discrete Algorithms*. ACM, New York, pp. 351–362.

RECEIVED APRIL 1997; REVISED JUNE 1998; ACCEPTED JULY 1998