

Wait-Free Synchronization

MAURICE HERLIHY

Digital Equipment Corporation

A *wait-free* implementation of a concurrent data object is one that guarantees that any process can complete any operation in a finite number of steps, regardless of the execution speeds of the other processes. The problem of constructing a wait-free implementation of one data object from another lies at the heart of much recent work in concurrent algorithms, concurrent data structures, and multiprocessor architectures. First, we introduce a simple and general technique, based on reduction to a consensus protocol, for proving statements of the form, “there is no wait-free implementation of X by Y .” We derive a hierarchy of objects such that no object at one level has a wait-free implementation in terms of objects at lower levels. In particular, we show that atomic read/write registers, which have been the focus of much recent attention, are at the bottom of the hierarchy: they cannot be used to construct wait-free implementations of many simple and familiar data types. Moreover, classical synchronization primitives such as *test&set* and *fetch&add*, while more powerful than *read* and *write*, are also computationally weak, as are the standard message-passing primitives. Second, nevertheless, we show that there do exist simple universal objects from which one can construct a wait-free implementation of any sequential object.

Categories and Subject Descriptors: D.1.3 [Programming Techniques]: Concurrent Programming; D.3.3 [Programming Languages]: Language Constructs—*abstract data types, concurrent programming structures*; D.4.1 [Operating Systems]: Process Management—*concurrency, mutual exclusion, synchronization*

General Terms: Algorithms, Languages, Verification

Additional Key Words and Phrases: Linearizability, wait-free synchronization

1. INTRODUCTION

A *concurrent object* is a data structure shared by concurrent processes. Algorithms for implementing concurrent objects lie at the heart of many important problems in concurrent systems. The traditional approach to implementing such objects centers around the use of *critical sections*: only one process at a time is allowed to operate on the object. Nevertheless, critical sections are poorly suited for asynchronous, fault-tolerant systems: if a faulty process is halted or delayed in a critical section, nonfaulty processes will also be unable to progress. Even in a failure-free system, a process can encounter unexpected delay as a result of a

A preliminary version of this paper appeared in the *Proceedings of the 7th ACM SIGACT-SIGOPS Symposium on Principles of Distributed Computing* (Aug. 1988), pp. 276–290.

Author’s address: Digital Equipment Corporation Cambridge Research Laboratory, One Kendall Square, Cambridge, MA 02139.

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.

© 1991 ACM 0164-0925/91/0100-0124 \$01.50

ACM Transactions on Programming Languages and Systems, Vol. 11, No. 1, January 1991, Pages 124–149.

page fault or cache miss, by exhausting its scheduling quantum, or if it is swapped out. Similar problems arise in heterogeneous architectures, where some processors may be inherently faster than others, and some memory locations may be slower to access.

A *wait-free* implementation of a concurrent data object is one that guarantees that any process can complete any operation in a finite number of steps, regardless of the execution speeds on the other processes. The wait-free condition provides fault-tolerance: no process can be prevented from completing an operation by undetected halting failures of other processes, or by arbitrary variations in their speed. The fundamental problem of wait-free synchronization can be phrased as follows:

Given two concurrent objects X and Y , does there exist a wait-free implementation of X by Y ?

It is clear how to show that a wait-free implementation exists: one displays it. Most of the current literature takes this approach. Examples include “atomic” registers from nonatomic “safe” registers [18], complex atomic registers from simpler atomic registers [4, 5, 16, 23, 25, 26, 29, 31], read-modify-write operations from combining networks [11, 15], and typed objects such as queues or sets from simpler objects [14, 19, 20].

It is less clear how to show that such an implementation does *not* exist. In the first part of this paper, we propose a simple new technique for proving statements of the form “there is no wait-free implementation of X by Y .” We derive a hierarchy of objects such that no object at one level can implement any object at higher levels (see Figure 1). The basic idea is the following: each object has an associated *consensus number*, which is the maximum number of processes for which the object can solve a simple consensus problem. In a system of n or more concurrent processes, we show that it is impossible to construct a wait-free implementation of an object with consensus number n from an object with a lower consensus number.

These impossibility results do not by any means imply that wait-free synchronization is impossible or infeasible. In the second part of this paper, we show that there exist *universal* objects from which one can construct a wait-free implementation of any object. We give a simple test for universality, showing that an object is universal in a system of n processes if and only if it has a consensus number greater than or equal to n . In Figure 1, each object at level n is universal for a system of n processes. A machine architecture or programming language is computationally powerful enough to support arbitrary wait-free synchronization if and only if it provides a universal object as a primitive.

Most recent work on wait-free synchronization has focused on the construction of atomic read/write registers [4, 5, 16, 18, 23, 25, 26, 29, 31]. Our results address a basic question: what are these registers good for? Can they be used to construct wait-free implementations of more complex data structures? We show that atomic registers have few, if any, interesting applications in this area. From a set of atomic registers, we show that it is impossible to construct a wait-free implementation of (1) common data types such as sets, queues, stacks, priority queues, or lists, (2) most if not all the classical synchronization primitives, such as *test&set*,

| Consensus Number | Object |
|------------------|--|
| 1 | read/write registers |
| 2 | test&set, swap, fetch&add, queue, stack |
| ⋮ | ⋮ |
| $2n - 2$ | n -register assignment |
| ⋮ | ⋮ |
| ∞ | memory-to-memory move and swap, augmented queue, compare&swap, fetch&cons, sticky byte |

Fig. 1. Impossibility and universality hierarchy.

compare&swap, and *fetch&add*, and (3) such simple memory-to-memory operations as *move* or memory-to-memory *swap*. These results suggest that further progress in understanding wait-free synchronization requires turning our attention from the conventional *read* and *write* operations to more fundamental primitives.

Our results also illustrate inherent limitations of certain multiprocessor architectures. The NYU ultracomputer project [10] has investigated architectural support for wait-free implementations of common synchronization primitives. They use combining networks to implement *fetch&add*, a generalization of *test&set*. IBM's RP3 [8] project is investigating a similar approach. The *fetch&add* operation is quite flexible: it can be used for semaphores, for highly concurrent queues, and even for database synchronization [11, 14, 30]. Nevertheless, we show that it is not universal, disproving a conjecture of Gottlieb et al. [11]. We also show that message-passing architectures such as hypercubes [28] are not universal either.

This paper is organized as follows. Section 2 defines a model of computation, Section 3 presents impossibility results, Section 4 describes some universal objects, and Section 5 concludes with a summary.

2. THE MODEL

Informally, our model of computation consists of a collection of sequential threads of control called *processes* that communicate through shared data structures called *objects*. Each object has a *type*, which defines a set of possible *states* and a set of primitive *operations* that provide the only means to manipulate that object. Each process applies a sequence of operations to objects, issuing an invocation and receiving the associated response. The basic correctness condition for concurrent systems is *linearizability* [14]: although operations of concurrent processes may overlap, each operation appears to take effect instantaneously at some point between its invocation and response. In particular, operations that do not overlap take effect in their “real-time” order.

2.1 I/O Automata

Formally, we model objects and processes using a simplified form of I/O automata [22]. Because the wait-free condition does not require any fairness or liveness conditions, and because we consider only finite sets of processes and objects, we

do not make use of the full power of the I/O automata formalism. Nevertheless, simplified I/O automata provide a convenient way to describe the basic structure of our model and to give the basic definition of what it means for one object to implement another. For brevity, our later constructions and impossibility results are expressed less formally using pseudocode. It is a straightforward exercise to translate this notation into I/O automata.

An *I/O automaton* A is a nondeterministic automaton with the following components¹:

- (1) $States(A)$ is a finite or infinite set of states, including a distinguished set of starting states.
- (2) $In(A)$ is a set of input events,
- (3) $Out(A)$ is a set of output events,
- (4) $Int(A)$ is a set of internal events,
- (5) $Step(A)$ is a transition relation given by a set of triples (s', e, s) , where s and s' are states and e is an event. Such a triple is called a *step*, and it means that an automaton in state s' can undergo a transition to state s , and that transition is associated with the event e .

If (s', e, s) is a step, we say that e is *enabled* in s' . I/O automata must satisfy the additional condition that inputs cannot be disabled; for each input event e and each state s' , there exist a state s and a step (s', e, s) .

An *execution fragment* of an automaton A is a finite sequence $s_0, e_1, s_1, \dots, e_n, s_n$ or infinite sequence s_0, e_1, s_1, \dots of alternating states and events such that each (s_i, e_{i+1}, s_{i+1}) is a step of A . An *execution* is an execution fragment where s_0 is a starting state. A *history fragment* of an automaton is the subsequence of events occurring in an execution fragment, and a *history* is the subsequence occurring in an execution.

A new I/O automaton can be constructed by composing a set of compatible I/O automata. (In this paper we consider only finite compositions.) A set of automata are compatible if they share no output or internal events. A state of the composed automaton S is a tuple of component states, and a starting state is a tuple of component starting states. The set of events of S , $Events(S)$, is the union of the components' sets of events. The set of output events of S , $Out(S)$, is the union of the components' sets of output events; the set of internal events, $Int(S)$, is the union of the components' sets of internal events; and the set of input events of S , $In(S)$, is $In(S) - Out(S)$, all the input events of S that are not output events for some component. A triple (s', e, s) is in $Steps(S)$ if and only if, for all component automata A , one of the following holds: (1) e is an event of A , and the projection of the step onto A is a step of A , or (2) e is not an event of A , and A 's state components are identical in s' and s . Note that composition is associative. If H is a history of a composite automaton and A a component automaton, $H|A$ denotes the subhistory of H consisting of events of A .

¹To remain consistent with the terminology of [14] we use "event" where Lynch and Tuttle use "operation," and "history" where they use "schedule."

2.2 Concurrent Systems

A *concurrent system* is a set of processes and a set of objects. Processes represent sequential threads of control, and objects represent data structures shared by processes. A process P is an I/O automaton with output events $\text{INVOKE}(P, op, X)$, where op is an operation² of object X , and input events $\text{RESPOND}(P, res, X)$, where res is a result value. We refer to these events as *invocations* and *responses*. Two invocations and responses match if their process and object names agree. To capture the notion that a process represents a single thread of control, we say that a process history is *well formed* if it begins with an invocation and alternates matching invocations and responses. An invocation is *pending* if it is not followed by a matching response. An *object* X has input events $\text{INVOKE}(P, op, X)$, where P is a process and op is an operation of the object, and output events $\text{RESPOND}(P, res, X)$, where res is a result value. Process and object names are unique, ensuring that process and object automata are compatible.

A *concurrent system* $\{P_1, \dots, P_n; A_1, \dots, A_m\}$ is an I/O automaton composed from processes P_1, \dots, P_n and objects A_1, \dots, A_m , where processes and objects are composed by identifying corresponding INVOKE and RESPOND events. A history of a concurrent system is well formed if each $H \upharpoonright P_i$ is well formed, and a concurrent system is well formed if each of its histories is well-formed. Henceforth, we restrict our attention to well-formed concurrent systems.

An execution is *sequential* if its first event is an invocation, and it alternates matching invocations and responses. A history is sequential if it is derived from a sequential execution. (Notice that a sequential execution permits process steps to be interleaved, but at the granularity of complete operations.) If we restrict our attention to sequential histories, then the behavior of an object can be specified in a particularly simple way: by giving pre- and postconditions for each operation. We refer to such a specification as a *sequential specification*. In this paper, we consider only objects whose sequential specifications are *total*: if the object has a pending invocation, then it has a matching enabled response. For example, a partial *deq* might be undefined when applied to an empty queue, while a total *deq* would return an exception. We restrict our attention to objects whose operations are total because it is unclear how to interpret the wait-free condition for partial operations. For example, the most natural way to define the effects of a partial *deq* in a concurrent system is to have it wait until the queue becomes nonempty, a specification that clearly does not admit a wait-free implementation.

Each history H induces a partial “real-time” order $<_H$ on its operations: $op_0 <_H op_1$ if the response for op_0 precedes the invocation for op_1 . Operations unrelated by $<_H$ are said to be *concurrent*. If H is sequential, $<_H$ is a total order. A concurrent system $\{P_1, \dots, P_n; A_1, \dots, A_m\}$ is *linearizable* if, for each history H , there exists a sequential history S such that

- (1) for all P_i , $H \upharpoonright P_i = S \upharpoonright P_i$;
- (2) $<_H \subset <_S$.

In other words, the history “appears” sequential to each individual process, and this apparent sequential interleaving respects the real-time precedence ordering

² Op may also include argument values.

of operations. Equivalently, each operation appears to take effect instantaneously at some point between its invocation and its response. A concurrent object A is linearizable [14] if, for every history H of every concurrent system $\{P_1, \dots, P_n; A_1, \dots, A_j, \dots, A_m\}$, $H \mid A_j$ is linearizable. A linearizable object is thus “equivalent” to a sequential object, and its operations can also be specified by simple pre- and postconditions. Henceforth, all objects are assumed to be linearizable. Unlike related correctness conditions such as sequential consistency [17] or strict serializability [24], linearizability is a *local* property: a concurrent system is linearizable if and only if each individual object is linearizable [14]. We restrict our attention to linearizable concurrent systems.

2.3 Implementations

An implementation of an object A is a concurrent system $\{F_1, \dots, F_n; R\}$, where the F_i are called *front-ends*, and R is called the *representation* object. Informally, R is the data structure that implements A , and F_i is the procedure called by process P_i to execute an operation. An object implementation is shown schematically in Figure 2.

(1) The external events of the implementation are just the external events of A : each input event $\text{INVOKE}(P_i, op, A)$ of A is an input event of F_i , and each output event $\text{RESPOND}(P_i, res, A)$ of A is an output event of F_i .

(2) The implementation has the following internal events: each input event $\text{INVOKE}(F_i, op, R)$ of R is composed with the matching output event of F_i , and each output event $\text{RESPOND}(F_i, res, R)$ of R is composed with the matching input event of F_i .

(3) To rule out certain trivial solutions, front-ends share no events; they communicate indirectly through R .

Let I_j be an implementation of A_j . I_j is *correct*, if for every history H of every system $\{P_1, \dots, P_n; A_1, \dots, I_j, \dots, A_m\}$, there exists a history H' of $\{P_1, \dots, P_n; A_1, \dots, A_j, \dots, A_m\}$, such that $H \mid \{P_1, \dots, P_n\} = H' \mid \{P_1, \dots, P_n\}$.

An implementation is wait-free if the following are true:

- (1) It has no history in which an invocation of P_i remains pending across an infinite number of steps of F_i .
- (2) If P_i has a pending invocation in a state s , then there exists a history fragment starting from s , consisting entirely of events of F_i and R , that includes the response to that invocation.

The first condition rules out unbounded busy-waiting: a front-end cannot take an infinite number of steps without responding to an invocation. The second condition rules out conditional waiting: F_i cannot block waiting for another process to make a condition true. Note that we have not found it necessary to make fairness or liveness assumptions: a wait-free implementation guarantees only that if R eventually responds to all invocations of F_i , then F_i will eventually respond to all invocations of P_i , independently of process speeds.

An implementation is *bounded wait-free* if there exists N such that there is no history in which an invocation of P_i remains pending across N steps of F_i .

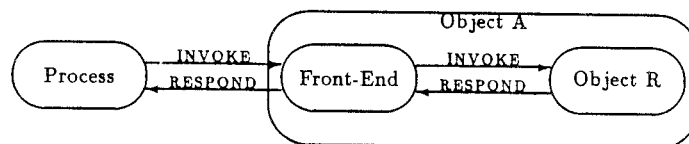


Fig. 2. Schematic view of object implementation.

Bounded wait-free implies wait-free, but not vice-versa. We use the wait-free condition for impossibility results and the bounded wait-free condition for universal constructions.

For brevity, we say that R implements A if there exists a wait-free implementation $\{F_1, \dots, F_n; R\}$ of A . It is immediate from the definitions that *implements* is a reflexive partial order on the universe of objects. In the rest of the paper, we investigate the mathematical structure of the *implements* relation. In the next section, we introduce a simple technique for proving that one object does *not* implement another, and in the following section we display some “universal” objects capable of implementing any other object.

3. IMPOSSIBILITY RESULTS

Informally, a *consensus protocol* is a system of n processes that communicate through a set of shared objects $\{X_1, \dots, X_m\}$. The processes each start with an input value from some domain \mathcal{S} ; they communicate with one another by applying operations to the shared objects; and they eventually agree on a common input value and halt. A consensus protocol is required to be

- (1) consistent: distinct processes never decide on distinct values;
- (2) wait-free: each process decides after a finite number of steps;
- (3) valid: the common decision value is the input to some process.

For our purposes, it is convenient to express the consensus problem using the terminology of abstract data types. A *consensus object* provides a single operation:

decide(input: value) returns(value)

A protocol’s sequential specification is simple: all *decide* operations return the argument value of the first *decide* (cf., Plotkin’s “sticky-bit” [27]). This common value is called the history’s decision value. A wait-free linearizable implementation of a consensus object is called a consensus protocol (cf., Fisher et al. [9]).

We investigate the circumstances under which it is possible to construct consensus protocols from particular objects. Most of the constructions presented in this paper use multireader/multiwriter registers in addition to the object of interest. For brevity we say “ X solves n -process consensus” if there exists a consensus protocol $\{F_1, \dots, F_n; W, X\}$, where W is a set of read/write registers, and W and X may be initialized to any state.

Definition 1. The consensus number for X is the largest n for which X solves n -process consensus. If no largest n exists, the consensus number is said to be infinite.

It is an immediate consequence of our definitions that if Y implements X , and X solves n -process consensus, then Y also solves n -process consensus.

THEOREM 1. *If X has consensus number n , and Y has consensus number $m < n$, then there exists no wait-free implementation of X by Y in a system of more than m processes.*

PROOF. As noted above, all front-end and object automata are compatible by definition, and thus their composition is well defined. Let $\{F_1, \dots, F_k; W, X\}$ be a consensus protocol, where $k > m$ and W is a set of read/write registers. Let $\{F'_1, \dots, F'_k; Y\}$ be an implementation of X . It is easily checked that $\{F_1, \dots, F_n; W, \{F'_1, \dots, F'_n; Y\}\}$ is wait-free, and because composition is associative, it is identical to $\{F_1 \cdot F'_1, \dots, F_n \cdot F'_n; W, Y\}$, where $F_i \cdot F'_i$ is the composition of F_i and F'_i . Since the former is a consensus protocol, so is the latter, contradicting the hypothesis that Y has consensus number m . \square

In the rest of this section, we consider a number of objects, displaying consensus protocols for some and impossibility results for others. For impossibility proofs, we usually assume the existence of a consensus protocol, and then derive a contradiction by constructing a sequential execution that forces the protocol to run forever. When constructing a consensus protocol for a particular linearizable object, we observe that the linearizability condition implies that if there exists an execution in which consensus fails, either because it is inconsistent, invalid, or it runs forever, then there exists an equivalent sequential execution with the same property. As a consequence, a consensus protocol is correct if and only if all its sequential executions are correct. For brevity, protocols are defined informally by pseudocode; their translations into I/O automata should be self-evident.

3.1 Atomic Read/Write Registers

In this section we show there exists no two-process consensus protocol using multireader/multiwriter atomic registers. First, we give some terminology. A protocol state is *bivalent* if either decision value is still possible; that is, the current execution can be extended to yield different decision values. Otherwise it is *univalent*. An *x-valent* state is a univalent state with eventual decision value x . A *decision step* is an operation that carries a protocol from a bivalent to a univalent state.

THEOREM 2. *Read/write registers have consensus number 1.*

PROOF. Assume there exists a two-process consensus protocol implemented from atomic read/write registers. We derive a contradiction by constructing an infinite sequential execution that keeps any such protocol in a bivalent state. If the processes have different input values, the validity condition implies that the initial state is bivalent. Consider the following sequential execution, starting from the initial state. In the first stage, P executes a sequence of operations (i.e., alternates matching invocation and response events) until it reaches a state where the next operation will leave the protocol in a univalent state. P must eventually reach such a state, since it cannot run forever, and it cannot block. In

Fig. 3. P reads first.

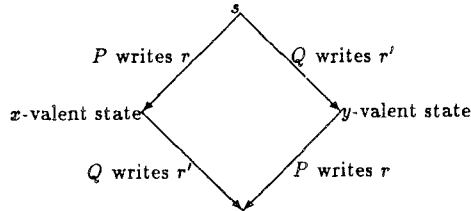
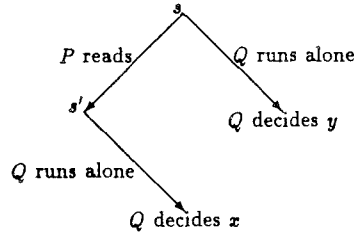


Fig. 4. P and Q write different registers.

the second stage, Q executes a sequence of operations until it reaches a similar state, and in successive stages, P and Q alternate sequences of operations until each is about to make a decision step. Because the protocol cannot run forever, it must eventually reach a bivalent state s in which any subsequent operation of either process is a decision step. Suppose P 's operation carries the protocol to an x -valent state, and Q 's operation carries the protocol to a y -valent state, where x and y are distinct.

(1) Suppose the decision step for one process, say P , is to read a shared register (Figure 3). Let s' be the protocol state immediately following the read. The protocol has a history fragment starting from s , consisting entirely of operations of Q , yielding decision value y . Since the states s and s' differ only in the internal state of P , the protocol has the same history fragment starting in s' , an impossibility because s' is x -valent.

(2) Suppose the processes write to different registers (Figure 4). The state that results if P 's write is immediately followed by Q 's is identical to the state that results if the writes occur in the opposite order, which is impossible, since one state is x -valent and the other is y -valent.

(3) Suppose the processes write to the same register (Figure 5). Let s' be the x -valent state immediately after P 's write. There exists a history fragment starting from s' consisting entirely of operations of P that yields the decision value x . Let s'' be the y -valent state reached if Q 's write is immediately followed by P 's. Because P overwrites the value written by Q , s' and s'' differ only in the internal states of Q , and therefore the protocol has the same history fragment starting from s'' , an impossibility since s'' is y valent. \square

Similar results have been shown by Loui and Abu-Amara [21], Chor et al. [6], and Anderson and Gouda [1]. Our contribution lies in the following corollary:

COROLLARY 1. *It is impossible to construct a wait-free implementation of any object with consensus number greater than 1 using atomic read/write registers.*

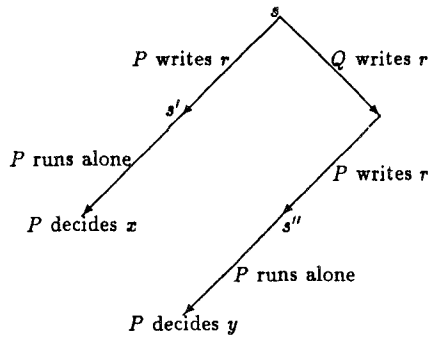
Fig. 5. P and Q write the same register.

Fig. 6. Read-modify-write.

```

RMW( $r$ : register,  $f$ : function) returns(value)
previous :=  $r$ 
 $r := f(r)$ 
return previous
end RMW

```

3.2 Read-Modify-Write Operations

Kruskal et al. [15] have observed that many, if not all, of the classical synchronization primitives can be expressed as *read-modify-write* operations, defined as follows. Let r be a register, and f a function from values to values. The operation $RMW(r, f)$ is informally defined by the procedure shown in Figure 6, which is executed atomically. If f is the identity, $RMW(r, f)$ is simply a *read* operation. A read-modify-write operation is *nontrivial* if f is not the identity function. Examples of well-known nontrivial read-modify-write operations include *test&set*, *swap*, *compare&swap*, and *fetch&add*. Numerous others are given in Kruskal et al. [15].

THEOREM 3. *A register with any nontrivial read-modify-write operation has a consensus number at least 2.*

PROOF. Since f is not the identity, there exists a value v such that $v \neq f(v)$. Let P and Q be two processes that share a two-register array *prefer*, where each entry is initialized to \perp , and a read-modify-write register r , initialized to v . P executes the protocol shown in Figure 7 (Q 's protocol is symmetric.)

Expressed in terms of the I/O automaton model, the read-modify-write register r is the object X , the *prefer* array is the set of atomic registers W , and the pseudocode in Figure 7 defines the front-end automaton for P . The front-end has three output events: the *write* and *RMW* invocations sent to r and *prefer*, and the decision value returned to P . Similarly, its input events are P 's invocation of *decide*, and the responses to the *write* and *RMW* invocations.

As noted above, because r and *prefer* are linearizable, it suffices to check correctness for sequential executions. The only operations that do not commute are the two read-modify-write operations applied to r . The protocol chooses P 's input if P 's operation occurs first, and Q 's input otherwise. \square

```

decide(input: value) returns(value)
  prefer[P] := input
  if RMW(r,f) = v
    then return prefer[P]
    else return prefer[Q]
  end if
end decide

```

Fig. 7. Read-modify-write: Two-process consensus.

COROLLARY 2. *It is impossible to construct a wait-free implementation of any nontrivial read-modify-write operation from a set of atomic read/write registers in a system with two or more processes.*

Although read-modify-write registers are more powerful than read/write registers, many common read-modify-write operations are still computationally weak. In particular, one cannot construct a wait-free solution to three process consensus using registers that support any combination of *read*, *write*, *test&set*, *swap*, and *fetch&add* operations. Let F be a set of functions indexed by an arbitrary set S . Define F to be *interfering* if for all values v and all i and j in S , either (1) f_i and f_j commute: $f_i(f_j(v)) = f_j(f_i(v))$, or (2) one function “overwrites” the other: either $f_i(f_j(v)) = f_i(v)$ or $f_j(f_i(v)) = f_j(v)$.

THEOREM 4. *There is no wait-free solution to three-process consensus using any combination of read-modify-write operations that apply functions from an interfering set F .*

PROOF. By contradiction. Let the three processes be P , Q , and R . As in the proof of Theorem 2, we construct a sequential execution leaving the protocol in bivalent state where every operation enabled for P and Q is a decision step, some operation of P carries the protocol to an x -valent state, and some operation of Q carries the protocol to a y -valent state, where x and y are distinct. By the usual commutativity argument, P and Q must operate on the same register; say, P executes $RMW(r, f_i)$ and Q executes $RMW(r, f_j)$.

Let v be the current value of register r . There are two cases to consider. First, suppose $f_i(f_j(v)) = f_j(f_i(v))$. The state s that results if P executes $RMW(r, f_i)$ and Q executes $RMW(r, f_j)$ is x -valent; thus there exists some history fragment consisting entirely of operations of R that yields decision value x . Let s' be the state that results if P and Q execute their operations in the reverse order. Since the register values are identical in s and s' , the protocol has the same history fragment starting in s' , contradicting the hypothesis that s' is y -valent.

Second, suppose $f_i(f_j(v)) = f_j(v)$. The state s that results if P executes $RMW(r, f_i)$ and Q executes $RMW(r, f_j)$ is x -valent; thus there exists some history fragment consisting entirely of operations of R that yields decision value x . Let s' be the state that results if Q alone executes its operation. Since the register values are identical in s and s' , the protocol has the same history fragment starting in s' , contradicting the hypothesis that s' is y -valent. \square

It follows that one cannot use any combination of these classical primitives to construct a wait-free implementation of any object with consensus number greater than 2.

```

compare&swap(r: register, old: value, new: value)
  returns(value)
  previous := r
  if previous = old
    then r := new
  end if
  return previous
end compare&swap

```

Fig. 8. Compare&Swap.

Fig. 9. Compare&Swap: n -process consensus.

```

decide(input: value) returns(value)
  first := compare&swap(r,  $\perp$ , input)
  if first = bottom
    then return input
  else return first
  end if
end decide

```

Another classical primitive is *compare&swap*, shown in Figure 8. This primitive takes two values: *old* and *new*. If the register's current value is equal to *old*, it is replaced by *new*; otherwise is left unchanged. The register's old value is returned.

THEOREM 5. *A compare&swap register has infinite consensus number.*

PROOF. In the protocol shown in Figure 9, the processes share a register r initialized to \perp . Each process attempts to replace \perp with its input; the decision value is established by the process that succeeds.

This protocol is clearly wait-free, since it contains no loops. Consistency follows from the following observations: (1) $r \neq \perp$ is a postcondition of *compare&swap*, and (2) for any $v \neq \perp$, the assertion $r = v$ is *stable*—once it becomes true, it remains true. Validity follows from the observation that if $r \neq \perp$, then r contains some process's input. \square

COROLLARY 3. *It is impossible to construct a wait-free implementation of a compare&swap register from a set of registers that support any combination of read, write, test&set, swap, or fetch&add operations in a system of three or more processes.*

3.3 Queues, Stacks, Lists, etc.

Consider a FIFO queue with two operations: *enq* places an item at the end of the queue, and *deq* removes the item from the head of the queue, returning an error value if the queue is empty.

THEOREM 6. *The FIFO queue has consensus number at least 2.*

PROOF. Figure 10 shows a two-process consensus protocol. The queue is initialized by enqueueing the value 0 followed by the value 1. As above, the processes share a two-element array *prefer*. P executes the protocol shown in Figure 10 (Q 's protocol is symmetric). Each process dequeues an item from the queue, returning its own preference if it dequeues 0, and the other's preference if it dequeues 1.

```

decide(input: value) returns(value)
  prefer[P] := input
  if deq(q) = 0
    then return prefer[P]
    else return prefer[Q]
  end if
end decide

```

Fig. 10. FIFO queues: Two-process consensus.

The protocol is wait-free, since it contains no loops. If each process returns its own input, then they must both have dequeued 0, violating the queue specification. If each returns the others' input, then they must both have dequeued 1, also violating the queue specification. Let the "winner" be the process that dequeues 0. Validity follows by observing that the winner's position in *prefer* is initialized before the first queue operation. \square

Trivial variations of this program yield protocols for stacks, priority queues, lists, sets, or any object with operations that return different results if applied in different orders.

COROLLARY 4. *It is impossible to construct a wait-free implementation of a queue, stack, priority queue, set, or list from a set of atomic read/write registers.*

Although FIFO queues solve two-process consensus, they cannot solve three-process consensus.

THEOREM 7. *FIFO queues have consensus number 2.*

PROOF. By contradiction. Assume we have a consensus protocol for processes P , Q , and R . As before, we maneuver the protocol to a state where P and Q are each about to make a decision step. Assume that P 's operation would carry the protocol to an x -valent state and Q 's to a y -valent state. The rest is a case analysis.

First, suppose P and Q both execute *deq* operations. Let s be the protocol state if P dequeues and then Q dequeues, and let s' be the state if the dequeues occur in the opposite order. Since s is x -valent, there exists a history fragment from s , consisting entirely of operations of R , yielding decision value x . But s and s' differ only in the internal states of P and Q ; thus the protocol has the same history fragment from s' , a contradiction because s' is y -valent.

Second, suppose P does an *enq* and Q a *deq*. If the queue is nonempty, the contradiction is immediate because the two operations commute: R cannot observe the order in which they occurred. If the queue is empty, then the y -valent state reached if Q dequeues and then P enqueues is indistinguishable to R from the x -valent state reached if P alone enqueues.

Finally, suppose both P and Q do *enq* operations. Let s be the state at the end of the following execution:

- (1) P and Q enqueue items p and q in that order.
- (2) Run P until it dequeues p . (Since the only way to observe the queue's state is via the *deq* operation, P cannot decide before it observes one of p or q .)
- (3) Run Q until it dequeues q .

Let s' be the state after the following alternative execution:

- (1) Q and P enqueue items q and p in that order.
- (2) Run P until it dequeues q .
- (3) Run Q until it dequeues p .

Clearly, s is x -valent and s' is y -valent. Both of P 's executions are identical until it dequeues p or q . Since P is halted before it can modify any other objects, Q 's executions are also identical until it dequeues p or q . By a now-familiar argument, a contradiction arises because s and s' are indistinguishable to R . \square

Trivial variations of this argument can be applied to show that many similar data types, such as sets, stacks, double-ended queues, and priority queues, all have consensus number 2.

A *message-passing architecture* (e.g., a hypercube [28]) is a set of processors that communicate via shared FIFO queues. Theorem 7 implies that message-passing architectures cannot solve three-process consensus or implement any object that can. Dolev et al. [7] give a related result: point-to-point FIFO message channels cannot solve two-process consensus. That result does not imply Theorem 7, however, because a queue item, unlike a message, is not “addressed” to any particular process, and hence it can be dequeued by anyone.

3.4 An Augmented Queue

Let us augment the queue with one more operation: *peek* returns but does not remove the first item in the queue.

THEOREM 8. *The augmented queue has infinite consensus number.*

PROOF. In the protocol shown in Figure 11, the queue q is initialized to *empty*, and each process enqueues its own input. The decision value is the input of the process whose *enq* occurs first.

As usual, the protocol is wait-free, since it contains no loops. Consistency follows from the following observations: (1) “the queue is nonempty” is a postcondition of each *enq*, and hence a precondition for each *peek*, and (2) for any v , “ v is the first item in the queue” is stable. Validity follows from the observation that the first item in the queue is some process’s input. \square

COROLLARY 5. *It is impossible to construct a wait-free implementation of the augmented queue from a set of registers supporting any combination of read, write, test&set, swap, or fetch&add operations.*

COROLLARY 6. *It is impossible to construct a wait-free implementation of the augmented queue from a set of regular queues.*

The *fetch&cons* operation atomically threads an item onto the front of a linked list. By an argument virtually identical to the one given for Theorem 8, a linked list with *fetch&cons* has infinite consensus number.

```

decide(input: value) returns(value)
  enq(q, input)
  return peek(q)
end decide

```

Fig. 11. Augmented FIFO queue: n -process consensus.

3.5 Memory-to-Memory Operations

Consider a collection of atomic read/write registers having one additional operation: *move* atomically copies the value of one register to another.³ We use the expression “ $a \leftarrow b$ ” to move the contents of b to a .

THEOREM 9. *An array of registers with move has infinite consensus number.*

PROOF. An n -process consensus protocol appears in Figure 12. The processes share two arrays: $prefer[1..n]$ and $r[1..n, 1..2]$, where $r[P, 1]$ is initialized to 1 and $r[P, 2]$ to 0, for $1 \leq P \leq n$. The protocol is clearly wait-free, since all loops are bounded.

To show consistency, we use the following assertions:

$$\begin{aligned} \mathcal{P}(P) &\equiv r[P, 1] = 0 \wedge r[P, 2] = 0 \\ \mathcal{Q}(P) &\equiv r[P, 2] = 1 \\ \mathcal{S}(P) &\equiv \mathcal{P}(P) \vee \mathcal{Q}(P) \end{aligned}$$

It is easily checked that $\mathcal{P}(P)$, $\mathcal{Q}(P)$, and $\mathcal{S}(P)$ are stable for each P , that $\mathcal{P}(P)$ and $\mathcal{Q}(P)$ are mutually exclusive, that $\mathcal{S}(P)$ is true after P executes statement 2, and that $\mathcal{S}(i)$ is true after each execution of statement 4. We say that a process P has *stabilized* if $\mathcal{S}(P)$ holds.

We claim that if $\mathcal{P}(P)$ holds for some P , then $\mathcal{Q}(Q)$ holds for some $Q < P$, and that every process between Q and P has stabilized. Let P be the least process for which $\mathcal{P}(P)$ holds. Since $r[P, 1]$ and $r[P, 2]$ are both 0, some $Q < P$ must have assigned 0 to $r[P, 1]$ (statement 4) before P executed statement 2. Q , however, executes statement 2 before statement 4; hence $\mathcal{S}(Q)$ holds. Since $\mathcal{P}(Q)$ is false by hypothesis, $\mathcal{Q}(Q)$ must hold. Moreover, if Q has assigned to $r[P, 1]$, then it has assigned to every $r[P', 1]$ for $Q < P' < P$; thus each such P' has stabilized.

Define the termination assertion as follows:

$$\mathcal{T}(P) \equiv \mathcal{Q}(P) \wedge (\forall Q > P) \mathcal{P}(Q).$$

\mathcal{T} is stable and it holds for at most one process. When P finishes the first loop (statements 3–5), every process greater than or equal to P has stabilized. If any of them satisfies \mathcal{T} , we are done. Otherwise, there exists a largest $Q < P$ satisfying $\mathcal{Q}(Q)$, and all the processes between P and Q have stabilized, implying that $\mathcal{T}(Q)$ holds. When P 's protocol terminates, it chooses the input of the unique Q satisfying $\mathcal{T}(Q)$. Since the termination assertion is stable, all processes agree.

Validity follows because $prefer[P]$ must be initialized before $\mathcal{T}(P)$ can become true. \square

³ Memory-to-memory *move* should not be confused with assignment; the former copies values between two public registers, while the latter copies value between public and private registers.

```

decide(input: value) returns(value)
  prefer[P] := input                                1
  r[P,2] ← r[P,1]                                  2
  for i in P+1 .. n do                              3
    r[i, 1] := 0                                    4
  end for                                           5
  for i in n .. 1 do                                 6
    if r[i,2] = 1                                   7
      then return prefer[i]                         8
    end if                                          9
  end for                                          10
end decide

```

Fig. 12. Memory-to-memory move: n -process consensus.

THEOREM 10. *An array of registers with memory-to-memory swap⁴ has infinite consensus number.*

PROOF. The protocol is shown in Figure 13. The processes share an array of registers $a[1 \dots n]$ whose elements are initialized to 0 and a single register r , initialized to 1. The first process to swap 1 into a wins. The protocol is wait-free because the loop is bounded. To show consistency, consider the following assertions where “ $\exists!P$ ” means “there exists a unique P .”

$$r = 1 \vee (\exists!P)a[P] = 1, \quad r = 0.$$

The first assertion is invariant, and the second is stable and becomes true after the first *swap*. It follows that each process observes a unique, stable P such that $a[P] = 1$.

Validity follows because each process initializes its position in *prefer* before executing a *swap*. \square

COROLLARY 7. *It is impossible to construct a wait-free implementation of memory-to-memory move or swap from a set of registers that support any combination of read, write, test&set, swap, or fetch&add operations.*

COROLLARY 8. *It is impossible to construct a wait-free implementation of memory-to-memory move or swap from a set of FIFO queues.*

3.6 Multiple Assignment

The expression

$$r_1, \dots, r_m := v_1, \dots, v_m$$

atomically assigns each value v_i to each register r_i .

THEOREM 11. *Atomic m -register assignment has consensus number at least m .*

⁴The memory-to-memory *swap* should not be confused with the read-modify-write *swap*; the former exchanges the values of two public registers, while the latter exchanges the value of a public register with a processor’s private register.


```

decide(input: value) returns(value)
  prefer[P] := input
  swap(a[P],r)
  for Q in 1 .. n do
    if a[Q] = 1
      then return prefer[Q]
    end if
  end for
end decide

```

Fig. 13. Memory-to-memory swap: n -process consensus.

PROOF. The protocol uses m “single-writer” registers r_1, \dots, r_m , where P_i writes to register r_i , and $m(m-1)/2$ “multiwriter” registers r_{ij} , where $i > j$, where P_i and P_j both write to register r_{ij} . All registers are initialized to \perp . Each process atomically assigns its input value to m registers: its single-writer register and its $m-1$ multiwriter registers. The decision value of the protocol is the first value to be assigned.

After assigning to its registers, a process determines the relative ordering of the assignments for two processes P_i and P_j as follows:

- (1) Read r_{ij} . If the value is \perp , then neither assignment has occurred.
- (2) Otherwise, read r_i and r_j . If r_i 's value is \perp , then P_j precedes P_i , and similarly for r_j .
- (3) If neither r_i nor r_j is \perp , reread r_{ij} . If its value is equal to the value read from r_i , then P_j precedes P_i , else vice-versa.

By repeating this procedure, a process can determine the value written by the earliest assignment. \square

This result can be improved.

THEOREM 12. *Atomic m -register assignment has consensus number at least $2m-2$.*

PROOF. Consider the following two-phase protocol. Each process has two single-writer registers, one for each phase, and each pair of processes share a register. Divide the processes into two predefined groups of $m-1$. In the first phase, each group achieves consensus within itself using the protocol from Theorem 11. In the second phase, each process atomically assigns its group's value to its phase-two single-writer register and the $m-1$ multiwriter registers shared with processes in the other group. Using the ordering procedures described above, the process constructs a directed graph G with the property that there is an edge from P_j to P_k if P_j and P_k are in different groups and the former's assignment precedes the latter's. It then locates a *source* process having at least one outgoing edge but no incoming edges, and returns that process's value. At least one process have performed an assignment; thus G has edges. Let Q be the process whose assignment is first in the linearization order. Q is a source, and it has an outgoing edge to every process in the other group; thus no process in the other group is also a source. Therefore, all source processes belong to the same group. \square

This algorithm is optimal with respect to the number of processes.

THEOREM 13. *Atomic m -register assignment has consensus number exactly $2m - 2$.*

PROOF. We show that atomic m -register assignment cannot solve $2m - 1$ process consensus for $m > 1$. By the usual construction, we can maneuver the protocol into a bivalent state s in which any subsequent operation executed by any process is a decision step. We refer to the decision value forced by each process as its default.

We first show that each process must have a “single-writer” register that it alone writes to. Suppose not. Let P and Q be processes with distinct defaults x and y . Let s' be the state reached from s if P performs its assignment, Q performs its assignment, and the other processes perform theirs. Because P went first, s' is x -valent. By hypothesis, every register written by P has been overwritten by another process. Let s'' be the state reached from s if P halts without writing, but all other processes execute in the same order. Because Q wrote first, s'' is y -valent. There exists a history fragment from s' , consisting entirely of operations of Q , with decision value x . Because the values of the registers are identical in s' and s'' , the protocol has the same history fragment from s'' , a contradiction because s'' is y -valent.

We next show that if P and Q have distinct default values, then there must be some register written only by those two processes. Suppose not. Let s' be the state reached from s if P performs its assignment and Q performs its assignment, followed by all other processes' assignments. Let s'' be the state reached by the same sequence of operations, except that P and Q execute their assignments in the reverse order. Because s' is x -valent, there exists a history fragment from s' consisting of operations of P with decision value x . But because every register written by both P and Q has been overwritten by some other process, the register values are the same in both s and s' ; hence the protocol has the same history fragment from s'' , a contradiction.

It follows that if P has default value x , and there are k processes with different default values, then P must assign to $k + 1$ registers. If there are $2m - 1$ processes which do not all have the same default, then some process must disagree with at least m other processes, and that process must assign to $m + 1$ registers. \square

The last theorem shows that consensus is irreducible in the following sense: it is impossible to achieve consensus among $2n$ processes by combining protocols that achieve consensus among at most $2m < 2n$ processes. If it were possible, one could implement each individual $2m$ -process protocol using $m - 1$ -register assignment, yielding a $2n$ -process consensus protocol, contradicting Theorem 13.

3.7 Remarks

Fischer et al. [9] have shown that there exists no two-process consensus protocol using message channels that permit messages to be delayed and reordered. That result does not imply Theorem 2, however, because atomic read/write registers lack certain commutativity properties of asynchronous message buffers. (In particular, [9, Lemma 1] does not hold.)

Dolev et al. [7] give a thorough analysis of the circumstances under which consensus can be achieved by message-passing. They consider the effects of

32 combinations of parameters: synchronous versus asynchronous processors, synchronous versus asynchronous communication, FIFO versus non-FIFO message delivery, broadcast versus point-to-point transmission, and whether *send* and *receive* are distinct primitives. Expressed in their terminology, our model has asynchronous processes, synchronous communication, and distinct *send* and *receive* primitives. We model *send* and *receive* as operations on a shared message channel object; whether delivery is FIFO and whether broadcast is supported depends on the type of the channel. Some of their results translate directly into our model: it is impossible to achieve two-process consensus by communicating through a shared channel that supports either broadcast with unordered delivery, or point-to-point transmission with FIFO delivery. Broadcast with ordered delivery, however, does solve n -process consensus.

A *safe* read/write register [18] is one that behaves like an atomic read/write register as long as operations do not overlap. If a *read* overlaps a *write*, however, no guarantees are made about the value read. Since atomic registers implement safe registers, safe registers cannot solve two-process consensus, and hence the impossibility results we derive for atomic registers apply equally to safe registers. Similar remarks apply to atomic registers that restrict the number of readers or writers.

Loui and Abu-Amara [21] give a number of constructions and impossibility results for consensus protocols using shared read-modify-write registers, which they call “test&set” registers. Among other results, they show that n -process consensus for $n > 2$ cannot be solved by read-modify-write operations on single-bit registers.

Lamport [19] gives a queue implementation that permits one enqueueing process to execute concurrently with one dequeueing process. With minor changes, this implementation can be transformed into a wait-free implementation using atomic read/write registers. Theorem 2 implies that Lamport’s queue cannot be extended to permit concurrent *deq* operations without augmenting the *read* and *write* operations with more powerful primitives.

A concurrent object implementation is *nonblocking* if it guarantees that some process will complete an operation in a finite number of steps, regardless of the relative execution speeds of the processes. The nonblocking condition guarantees that the system as a whole will make progress despite individual halting failures or delays. A wait-free implementation is necessarily nonblocking, but not vice-versa, since a nonblocking implementation may permit individual processes to starve. The impossibility and universality results presented in this paper hold for nonblocking implementations as well as wait-free implementations.

Elsewhere [14], we give a nonblocking implementation of a FIFO queue, using *read*, *fetch&add*, and *swap* operations, which permits an arbitrary number of concurrent *enq* and *deq* operations. Corollary 5 implies that this queue implementation cannot be extended to support a nonblocking *peek* operation without introducing more powerful primitives.

4. UNIVERSALITY RESULTS

An object is *universal* if it implements any other object. In this section, we show that any object with consensus number n is universal in a system of n (or fewer)

processes. The basic idea is the following: we represent the object as a linked list, where the sequence of cells represents the sequence of operations applied to the object (and hence the object's sequence of states). A process executes an operation by threading a new cell on to the end of the list. When the cell becomes sufficiently old, it is reclaimed and reused. Our construction requires $O(n^3)$ memory cells to represent the object, and $O(n^3)$ worst case time to execute each operation. We assume cells can hold integers of unbounded size. Our presentation is intended to emphasize simplicity, and omits many obvious optimizations.

Let *INVOC* be the object's domain of invocations, *RESULT* its domain of results, and *STATE* its domain of states. An object's behavior may be specified by the following relation:

$$\text{apply} \subset \text{INVOC} \times \text{STATE} \times \text{STATE} \times \text{RESULT}.$$

This specification means that applying operation p in state s leaves the object in a state s' and returns result value r , where $\langle p, s, s', r \rangle \in \text{apply}$. *Apply* is a relation (rather than a function) because the operation may be nondeterministic. For brevity, we use the notation $\text{apply}(p, s)$ to denote an arbitrary pair $\langle s', r \rangle$ such that $\langle p, s, s', r' \rangle \in \text{apply}$.

4.1 The Algorithm

An object is represented by a doubly linked list of cells having the following fields:

- (1) *Seq* is the cell's sequence number in the list. This field is zero if the cell is initialized but not yet threaded onto the list, and otherwise it is positive. Sequence numbers for successive cells in the list increase by one.
- (2) *Inv* is the invocation (operation name and argument values).
- (3) *New* is a consensus object whose value is the pair $\langle \text{new.state}, \text{new.result} \rangle$. The first component is the object's state following the operation, and the second is the operation's result value, if any.
- (4) *Before* is a pointer to the previous cell in the list. This field is used only for free storage management.
- (5) *After* is a consensus object whose value is a pointer to the next cell in the list.

If c and d are cells, the function $\text{max}(c, d)$ returns the cell with the higher sequence number.

Initially, the object is represented by a unique *anchor* cell with sequence number 1, holding a creation operation and an initial state.

The processes share the following data structures:

- (1) *Announce* is an n -element array whose P th element is a pointer to the cell P is currently trying to thread onto the list. Initially all elements point to the anchor cell.
- (2) *Head* is an n -element array whose P th element is a pointer to the last cell in the list that P has observed. Initially all elements point to the anchor cell.

Let $\max(\text{head})$ be $\max(\text{head}[1].\text{seq}, \dots, \text{head}[n].\text{seq})$, and let “ $c \in \text{head}$ ” denote the assertion that a pointer to cell c has been assigned to $\text{head}[Q]$, for some Q .

We use the following auxiliary variables:

- (1) $\text{concur}(P)$ is the set of cells whose addresses have been stored in the head array since P 's last announcement;
- (2) $\text{start}(P)$ is the value of $\max(\text{head})$ at P 's last announcement. Notice that

$$|\text{concur}(P)| + \text{start}(P) = \max(\text{head}). \quad (1)$$

Auxiliary variables do not affect the protocol's control flow; they are present only to facilitate proofs.

The protocol for process P is shown in Figure 14. In this figure, “ $v: T := e$ ” declares and initializes variable v of type T to a value e , and the type “*cell” means “pointer to cell.” Sequences of statements enclosed in angle brackets are executed atomically. In each of these compound statements, only the first affects shared data or control flow; the remainder are “bookkeeping operations” that update auxiliary variables. For readability, auxiliary variables are shown in italics.

Informally, the protocol works as follows. P allocates and initializes a cell to represent the operation (statement 1). It stores a pointer to the cell in $\text{announce}[P]$ (statement 2), ensuring that if P itself does not succeed in threading its cell onto the list, some other process will. To locate a cell near the end of the list, P scans the head array, setting $\text{head}[P]$ to the cell with the maximal sequence number (statement 3). P then enters the main loop of the protocol (statement 4), which it executes until its own cell has been threaded onto the list (detected when its sequence number becomes nonzero). P chooses a process to “help” (statement 6), and checks whether that process has an unthreaded cell (statement 7). If so, then P will try to thread it; otherwise it tries to thread its own cell. (If this helping step were omitted, the protocol would be nonblocking rather than wait-free.) P tries to set $\text{head}[P].\text{after}$ to point to the cell it is trying to thread (statement 8). The after field must be a consensus cell to ensure that only one process succeeds in setting it. Whether or not P succeeds, it then initializes the remaining fields of the next cell in the list. Because the operation may be nondeterministic, different processes may try to set the new field to different values, so this field must be a consensus object (statement 9). The values of the other fields are computed deterministically, so they can simply be written as atomic registers (statements 10 and 11). For brevity, we say that a process *threads* a cell in statement 7 if the *decide* operation alters the value of the after field, and it *announces* a cell at statement 2 when it stores the cell's address in announce .

LEMMA 1. *The following assertion is invariant:*

$$|\text{concur}(P)| > n \Rightarrow \text{announce}(P) \in \text{head}.$$

PROOF. If $|\text{concur}(P)| > n$, then $\text{concur}(P)$ includes successive cells q and r with respective sequence numbers equal to $P - 1 \bmod n$ and $P \bmod n$, threaded by processes Q and R . Because q is in $\text{concur}(P)$, Q threads q after P 's announcement. Because R cannot modify an unthreaded cell, R reads

```

universal(what: INVOC) returns(RESULT)
  mine: cell := {seq: 0,                                1
    inv: what,
    new: create(consensus_object),
    before: create(consensus_object)
    after: null}
  ⟨announce[P] := mine; start(P) := max(head)⟩          2
  for each process Q do                                3
    head[P] := max(head[P], head[Q])
  end for
  while announce[P].seq = 0 do                          4
    c: *cell := head[P]                                 5
    help: *cell := announce[(c.seq mod n) + 1]          6
    if help.seq = 0                                    7
      then prefer := help
      else prefer := announce[P]
    end if
    d := decide(c.after, prefer)                        8
    decide(d.new, apply(d.inv, c.new.state))            9
    d.before := c                                     10
    d.seq := c.seq + 1                                 11
    ⟨head[P] := d; (∀Q) concur(Q) := concur(Q) ∪ {d}⟩ 12
  end while
  ⟨head[P] := announce[P]; (∀Q) concur(Q) := concur(Q) ∪ {d}⟩ 13
  return (announce[P].new.result)                     14
end universal

```

Fig. 14. A universal construction.

$announce[P]$ (statement 5) after Q threads q . It follows that R reads $announce[P]$ after P 's announcement, and therefore either $announce[P]$ is already threaded, or r is p . \square

Lemma 1 places a bound on the number of cells that can be threaded while an operation is in progress. We now give a sequence of lemmas showing that when P finishes scanning the $head$ array, either $announce[P]$ is threaded, or $head[P]$ lies within $n + 1$ cells of the end of the list.

LEMMA 2. *The following assertion is invariant:*

$$\max(head) \geq start(P).$$

PROOF. The sequence number for each $head[Q]$ is nondecreasing. \square

LEMMA 3. *The following is a loop invariant for statement 3:*

$$\max(head[P], head[Q], \dots, head[n]) \geq start(P).$$

where Q is the loop index.

PROOF. When Q is 1, the assertion is implied by Lemma 2. The truth of the assertion is preserved at each iteration, when $head[P]$ is replaced by $\max(head[P], head[Q])$. \square

LEMMA 4. *The following assertion holds just before statement 4:*

$$\text{head}[P].\text{seq} \geq \text{start}(P).$$

PROOF. After the loop at statement 3, $\max(\text{head}[P], \text{head}[Q], \dots, \text{head}[n])$ is just $\text{head}[P].\text{seq}$, and the result follows from Lemma 3. \square

LEMMA 5. *The following is invariant:*

$$|\text{concur}(P)| \geq \text{head}[P].\text{seq} - \text{start}(P) \geq 0.$$

PROOF. The lower bound follows from Lemma 4, and the upper bound follows from (1). \square

THEOREM 14. *The protocol in Figure 14 is correct and bounded wait-free.*

PROOF. Linearizability is immediate, since the order in which cells are threaded is clearly compatible with the natural partial order of the corresponding operations.

The protocol is bounded wait-free because P can execute the main loop no more than $n + 1$ times. At each iteration, $\text{head}[P].\text{seq}$ increases by one. After $n + 1$ iterations, Lemma 5 implies that

$$|\text{concur}(P)| \geq \text{head}[P].\text{seq} - \text{start}(P) \geq n.$$

Lemma 1 implies that $\text{announce}[P]$ must be threaded. \square

4.2 Memory Management

In this section we discuss how cells are allocated and reclaimed. To reclaim a cell, we assume each consensus object provides a *reset* operation that restores the object to a state where it can be reused for a new round of consensus. Our construction resets a consensus object only when there are no concurrent operations in progress.

The basic idea is the following: a process executing an operation will traverse no more than $n + 1$ cells before its cell is threaded (Theorem 14). Conversely, each cell will be traversed no more than $n + 1$ times. When a process is finished threading its cell, it *releases* each of the $n + 1$ preceding cells by setting a bit. When a cell has been released $n + 1$ times, it is safe to recycle it. Each cell holds an additional field, an array *released* of $n + 1$ bits, initially all *false*. When a process completes an operation, it scans the $n + 1$ earlier cells, setting *released*[i] to *true* in the cell at distance i .

Each process maintains a private pool of cells. When a process needs to allocate a new cell, it scans its pool, and reinitializes the first cell whose *released* bits are all *true*. We assume here that each object has its own pool; in particular, the cell's new sequence number exceeds its old sequence number. While a process P is allocating a new cell, the list representing an object includes at most $n - 1$ incomplete operations, and each such cell can inhibit the reclamation of at most $n + 1$ cells. To ensure that P will find a free cell, it needs a pool of at least n^2 cells. Note that locating a free cell requires at worst $O(n^3)$ *read* operations, since the process may have to scan n^2 cells, and each cell requires reading $n + 1$ bits. If an atomic *fetch&add* operation is available, then a counter can be used instead of the *released* bits, and a free cell can be located in $O(n^2)$ *read* operations.

The proof of Lemma 1 remains unchanged. For Lemma 2, we observe that a cell can be reclaimed only if it is followed in the list by at least $n + 1$ other cells; hence reclaiming a cell cannot affect the value of $\max(head)$. The statement of Lemma 4 needs to be strengthened:

LEMMA 6. *The following assertion holds just before statement 4:*

$$announce[P] \in head \vee head[P].seq \geq start(P).$$

PROOF. When P announces its cell, there is some process Q such that $head[Q]$ has sequence number greater than or equal to $start(P)$. This cell can be reclaimed only if $n + 1$ other cells are threaded in front of it, implying that $|concur(P)| \geq n + 1$, and hence that $announce[P] \in head$ (Lemma 1). \square

The proof of Theorem 14 proceeds as before. There is one last detail to check: whether P 's cell has not been threaded by the time it finishes scanning $head$; then we claim that none of the cells it traverses will be reclaimed while the operation is in progress. Lemma 1 states that the list cannot have grown by more than n cells since P 's announcement; thus every cell reachable from $head[P]$ lies within $n + 1$ cells of the end of the list, or of $announce[P]$ if it is threaded. In either case, those cells cannot be reclaimed while P 's operation is in progress, since they must have at least one *released* bit unset.

4.3 Remarks

The first universal construction [13] used unbounded memory. Plotkin [27] describes a universal construction employing “sticky-byte” registers, a kind of write-once memory. In Plotkin's construction, cells are allocated from a common pool and reclaimed in a way similar to ours. The author [12] describes a universal construction using *compare&swap* that is currently being implemented on a multiprocessor.

A *randomized wait-free* implementation of a concurrent object is one that guarantees that any process can complete any operation in a finite *expected* number of steps. Elsewhere [2], we give a randomized consensus protocol using atomic registers whose expected running time is polynomial in the number of processes. This protocol has several important implications. If the wait-free guarantee is allowed to be probabilistic in nature, then the hierarchy shown in Figure 1 collapses because atomic registers become universal. Moreover, combining the randomized consensus protocol with our universal construction yields a polynomial-time randomized universal construction. Bar-Noy and Dolev [3] have adapted our randomized consensus protocol to a message-passing model; that protocol can be used to manage randomized wait-free *replicated* data objects.

5. CONCLUSIONS

Wait-free synchronization represents a qualitative break with the traditional locking-based techniques for implementing concurrent objects. We have tried to suggest here that the resulting theory has a rich structure, yielding a number of unexpected results with consequences for algorithm design, multiprocessor architectures, and real-time systems. Nevertheless, many interesting problems remain unsolved. Little is known about lower bounds for universal constructions,

both in terms of time (rounds of consensus) and space (number of cells). The *implements* relation may have additional structure not shown in the impossibility hierarchy of Figure 1. For example, can atomic registers implement any object with consensus number 1 in a system of two or more processes? Can *fetch&add* implement any object with consensus number 2 in a system of three or more processes? Does the *implements* relation have a different structure for bounded wait-free, wait-free, or nonblocking synchronization? Finally, little is known about practical implementation techniques.

ACKNOWLEDGMENTS

The author is grateful to Jim Aspnes, Vladimir Lanin, Michael Merritt, Serge Plotkin, Mark Tuttle, and Jennifer Welch for many invaluable suggestions.

REFERENCES

1. ANDERSON, J. H., AND GOUDA, M. G. The virtue of patience: Concurrent programming with and without waiting. Private communication.
2. ASPNES, J., AND HERLIHY, M. P. Fast randomized consensus using shared memory. *J. Algorithm.* 11 (1990), 441–461.
3. BAR-NOY, A., AND DOLEV, D. Shared memory vs. message-passing in an asynchronous distributed environment. In *Proceedings of the 8th ACM SIGACT-SIGOPS Symposium on Principles of Distributed Computing* (Aug. 1989), pp. 307–318.
4. BLOOM, B. Constructing two-writer atomic registers. In *Proceedings of the 6th ACM Symposium on Principles of Distributed Computing* (1987), pp. 249–259.
5. BURNS, J. E., AND PETERSON, G. L. Constructing multi-reader atomic values from non-atomic values. In *Proceedings of the 6th ACM Symposium on Principles of Distributed Computing* (1987), pp. 222–231.
6. CHOR, B., ISRAELI, A., AND LI, M. On processor coordination using asynchronous hardware. In *Proceedings of the 6th ACM Symposium on Principles of Distributed Computing* (1987), pp. 86–97.
7. DOLEV, D., DWORK, C., AND STOCKMEYER, L. On the minimal synchronism needed for distributed consensus. *J. ACM* 34, 1 (Jan. 1987), 77–97.
8. PFISTER, G. H., ET AL. The IBM research parallel processor prototype (rp3): Introduction and architecture. Presented at the International Conference on Parallel Processing, 1985.
9. FISCHER, M. J., LYNCH, N. A., AND PATERSON, M. S. Impossibility of distributed consensus with one faulty process. *J. ACM* 32, 2 (April 1985), 374–382.
10. GOTTLIEB, A., GRISHMAN, R., KRUSKAL, C. P., MCAULIFFE, K. P., RUDOLPH, L., AND SNIR, M. The NYU ultracomputer—Designing a mind parallel computer. *IEEE Trans. Comput. C-32*, 2 (Feb. 1984) 175–189.
11. GOTTLIEB, A., LUBACHEVSKY, B. D., AND RUDOLPH, L. Basic techniques for the efficient coordination of very large numbers of cooperating sequential processors. *ACM Trans. Program. Lang. Syst.* 5, 2 (April 1983), 164–189.
12. HERLIHY, M. P. A methodology for implementing highly concurrent data structures. In *Proceedings of the 2nd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming* (March 1990).
13. HERLIHY, M. P. Impossibility and universality results for wait-free synchronization. In *Proceedings of the 7th ACM SIGACT-SIGOPS Symposium on Principles of Distributed Computing* (Aug. 1988), pp. 276–290.
14. HERLIHY, M. P., AND WING, J. M. Axioms for concurrent objects. In *Proceedings of the 14th ACM Symposium on Principles of Programming Languages* (Jan. 1987), pp. 13–26.
15. KRUSKAL, C. P., RUDOLPH, L., AND SNIR, M. Efficient synchronization on multiprocessors with shared memory. In *Proceedings of the 5th ACM SIGACT-SIGOPS Symposium on Principles of Distributed Computing* (Aug. 1986).

16. LAMPORT, L. Concurrent reading and writing. *Commun. ACM* 20, 11 (Nov. 1977), 806–811.
17. LAMPORT, L. How to make a multiprocessor computer that correctly executes multiprocess programs. *IEEE Trans. Comput. C-28*, 9 (Sept. 1979), 690.
18. LAMPORT, L. On interprocess communication, parts i and ii. *Distributed Comput.* 1 (1986), 77–101.
19. LAMPORT, L. Specifying concurrent program modules. *ACM Trans. Program Lang. Syst.* 5, 2 (April 1983), 190–222.
20. LANIN, V., AND SHASHA, D. Concurrent set manipulation without locking. In *Proceedings of the 7th ACM Symposium on Principles of Database Systems* (March 1988), pp. 211–220.
21. LOUI, M. C., AND ABU-AMARA, H. H. *Memory Requirements for Agreement Among Unreliable Asynchronous Processes*, vol. 4. Jai Press, Greenwich, Conn., 1987, pp. 163–183.
22. LYNCH, N. A., AND TUTTLE, M. R. An introduction to input/output automata. Tech. Rep. MIT/LCS/TM-373, MIT Laboratory for Computer Science, Nov. 1988.
23. NEWMAN-WOLFE, R. A protocol for wait-free, atomic, multi-reader shared variables. In *Proceedings of the 6th ACM Symposium on Principles of Distributed Computing* (1987) pp. 232–249.
24. PAPADIMITRIOU, C. H. The serializability of concurrent database updates. *J. ACM* 26, 4 (Oct. 1979), 631–653.
25. PETERSON, G. L. Concurrent reading while writing. *ACM Trans. Program Lang. Syst.* 5, 1 (Jan. 1983), 46–55.
26. PETERSON, G. L., AND BURNS, J. E. Concurrent reading while writing II: The multi-writer case. Tech. Rep. GIT-ICS-86/26, Georgia Institute of Technology, Dec. 1986.
27. PLOTKIN, S. A. Sticky bits and universality of consensus. In *Proceedings of the 8th ACM Symposium on Principles of Distributed Computing* (1989), pp. 159–176.
28. SEITZ, C. L. The cosmic cube. *Commun. ACM* 28, 1 (Jan. 1985).
29. SINGH, A. K., ANDERSON, J. H., AND GOUDA, M. G. The elusive atomic register revisited. In *Proceedings of the 6th ACM Symposium on Principles of Distributed Computing*, (1987), pp. 206–221.
30. STONE, H. S. Database applications of the fetch-and-add instruction. *IEEE Trans. Comput. C-33*, 7 (July 1984), 604–612.
31. VITANYI, P., AND AWERBUCH, B. Atomic shared register access by asynchronous hardware. In *Proceedings of the 27th IEEE Symposium on Foundations of Computer Science*, (1986), pp. 223–243. See also *SIGACT News* 18, 4 (Summer, 1987), errata.

Received May 1989; revised February 1990, accepted March 1990