

# On the Inherent Weakness of Conditional Synchronization Primitives

Faith Fich<sup>\*</sup>  
Department of  
Computer Science  
University of Toronto  
fich@cs.toronto.edu

Danny Hendler<sup>†</sup>  
School of Computer Science  
Tel-Aviv University  
Tel Aviv, Israel 69978  
hendlerd@post.tau.ac.il

Nir Shavit  
Tel-Aviv University &  
Sun Microsystems  
Laboratories  
shanir@sun.com

## ABSTRACT

The “wait-free hierarchy” classifies multiprocessor synchronization primitives according to their power to solve consensus. The classification is based on assigning a number  $n$  to each synchronization primitive, where  $n$  is the maximal number of processes for which deterministic wait-free consensus can be solved using instances of the primitive and *read write* registers. Conditional synchronization primitives, such as *Compare-and-Swap* and *Load-Linked/Store-Conditional*, can implement deterministic wait-free consensus for any number of processes (they have consensus number  $\infty$ ), and are thus considered to be among the strongest synchronization primitives; *Compare-and-Swap* and *Load-Linked/Store-Conditional* have consequently become the synchronization primitives of choice, and have been implemented in hardware in many multiprocessor architectures.

This paper shows that, though they are strong in the context of consensus, conditional synchronization primitives are not efficient in terms of memory space for implementing many key objects. Our results hold for starvation-free implementations of mutual exclusion, and for wait-free implementations of a large class of concurrent objects, that we call *Visible(n)*. Roughly, *Visible(n)* is a class that includes all objects that support some operation that must perform a “visible” write before it terminates. *Visible(n)* includes many useful objects; some examples are: counters, stacks, queues, swap, fetch-and-add, and single-writer snapshot objects. We show that at least  $n$  conditional registers are required by any such implementation, even if registers are of unbounded size. We also obtain tradeoffs between time and space for  $n$ -process wait-free implementations of any *one-time* object in *Visible(n)*. All these results hold for both deterministic and randomized implementations.

<sup>\*</sup>Supported in part by grants from the Natural Sciences and Engineering Research Council of Canada, and from Sun Microsystems.

<sup>†</sup>Supported in part by a grant from Sun Microsystems.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

PODC'04, July 25–28, 2004, St. Johns, Newfoundland, Canada.  
Copyright 2004 ACM 1-58113-802-4/04/0007 ...\$5.00.

Starvation-free mutual exclusion and wait-free implementations of some objects in *Visible(n)* (e.g. counters, swap and fetch-and-add) can be implemented by  $O(1)$  non-conditional primitives. Thus we believe that basing multiprocessor synchronization solely on conditional synchronization primitives might not be the best design choice.

## Categories and Subject Descriptors

C.1.4.1 [Computer Systems Organization]: Processor Architectures—*Parallel Architectures, Distributed Architectures*

## General Terms

Theory

## Keywords

Lower bounds, wait-freedom, synchronization primitives, Compare-and-Swap, Load-Linked, Store-Conditional, Test-and-Set

## 1. INTRODUCTION

A *wait-free* implementation of a concurrent object is one that guarantees that any process can complete an operation in a finite number of its own steps. Wait-freedom implies *lock-freedom*, thus it also provides for better robustness, as processes coordinate without using locks, and so the inherent problems of locks, such as deadlock, convoying, and priority-inversion, are avoided. In 1991, Herlihy introduced his influential “wait-free hierarchy” [7], which classifies multiprocessor synchronization primitives according to their power to solve consensus. The classification is based on assigning a number  $n$  to each synchronization primitive, where  $n$  is the maximal number of processes for which deterministic wait-free consensus can be solved using instances of the primitive and read-write registers. The strongest synchronization primitives according to this classification are those that can solve consensus for *any* number of processes, and hence have a consensus number of  $\infty$ .

Synchronization primitives such as *Compare-and-Swap* (denoted *CAS*) and *Load-Linked/Store-Conditional* (denoted *LL/SC*) have consensus number  $\infty$  and are thus among the strongest primitives according to Herlihy’s hierarchy; to some extent because of that, they have become the synchronization primitives of choice, and have been implemented in

hardware in many multiprocessor architectures. (For example, *CAS* has been implemented in Motorola 680x0, IBM 370 and *SPARC*<sup>®</sup> architectures; *LL/SC* has been implemented in MIPS<sup>®</sup>, PowerPC<sup>®</sup> and DECAlpha<sup>®</sup> architectures.) *CAS* and *LL/SC* belong to a class of primitives which is called by the slightly misleading name *conditional synchronization primitives* (or simply *conditionals*). Conditionals are primitive operations that modify the value of the register  $r$  on which they operate only if its value, just before the operation starts, is a *specific* value,  $v_w$ , that depends on the input,  $w$ , of the conditional operation. As an example, if some process issues a *CAS*( $addr$ ,  $old$ ,  $new$ ) operation, the operation changes the contents of the register  $addr$  to  $new$  only if  $addr$ 's value just before the operation starts equals  $old$ ; otherwise, the operation is not visible to other processes (except that it might have an adverse effect on memory contention and process latency).

Several researchers have previously obtained results suggesting that conditionals are no stronger than read-write registers in the context of mutual exclusion. Cypher [4] has obtained a lower bound of  $\Omega(\log \log n / \log \log \log n)$  remote memory references for  $n$ -process mutual exclusion based on read-write registers and registers that support conditional operations. This bound was later improved by Anderson and Kim [1] to  $\Omega(\log n / \log \log n)$  remote memory references. In their survey of shared-memory mutual exclusion [2], Anderson and Kim comment on the "...unexpected weakness" of conditionals in general, and *CAS* in particular, which is "...still widely regarded as being the most useful of all primitives to provide in hardware". Our work shows that conditionals are relatively weak not just in the context of mutual exclusion time-complexity. Rather, we show that conditionals are inefficient in terms of memory space for implementing many widely used distributed objects. Collectively, these results imply that basing multiprocessor synchronization solely on conditional synchronization primitives might not be the best design choice.

In the article that introduced *covering* arguments [3], Burns and Lynch obtain a lower bound of  $n$  on the number of read-write registers required to solve mutual exclusion. In the context of mutual exclusion, our work extends their work for *starvation-free* mutual exclusion to obtain bounds for implementations that use conditionals. To do this, we use generalized covering arguments that can be applied to conditionals. Some of the ideas that we use in our proofs for long-lived objects are also similar to those used in their lower bound proof.

Fich, Herlihy, and Shavit, [6] obtain a lower bound of  $\Omega(\sqrt{n})$  on the space complexity of randomized wait-free implementations of  $n$ -process consensus that use *historyless* objects. These are objects, such as read-write registers, *test&set*, or *swap*, whose state depends only on the last non-trivial operation (for example, write) that was applied to them. Jayanti, Tan, and Toueg [10] show time and space lower bounds of  $n - 1$  for lock-free implementations of several objects (e.g. counters and  $k$ -valued *CAS*) from historyless objects and resettable consensus. Note that *test&set* is a historyless conditional.

A *long-lived* object is one whose operations may be called by a process any number of times. In contrast, the operations of *one-time* objects may be called by a process only once. For one-time objects, we show tradeoffs between memory space and the number of writes, and between memory

space and the number of *memory stalls* incurred by high-level operations. The concept of memory stalls was defined by Dwork, Herlihy, and Waarts, in a paper [5] that introduced a formal model to capture the phenomenon of memory contention in shared memory multiprocessors.

## 1.1 Our Results

Our results apply to implementations of any object that supports some operation that must perform a "visible" write before it terminates. This is a large class of concurrent objects that includes well-known objects such as counters (both linearizable and non-linearizable), stacks, queues, swap, fetch-and-add, and single-writer snapshot. For all such long lived objects, we show the following.

- Any  $n$ -process wait-free implementation that uses either only read-write registers or registers that support only conditional operations, requires at least  $n$  registers;
- Any  $n$ -process wait-free implementation that uses registers on which both read-write *and* conditional operations can be applied, requires at least  $\frac{n}{2}$  such registers.

For all such one-time objects, we show the following tradeoffs for any wait-free implementation that uses  $m < n$  registers that support only read, write and/or conditional operations.

- Either  $m > \sqrt{n}$ , or the amortized number of writes performed by high-level operations is in  $\Omega(\sqrt{n})$ ;
- Either  $m > n^{\frac{2}{3}}$ , or the amortized number of memory stalls incurred by high-level operations is in  $\Omega(n^{\frac{2}{3}})$ .

All of the above results apply also to starvation-free implementations of mutual exclusion. These results apply to both deterministic and randomized implementations, even if they use registers of unbounded size.

The rest of the paper is organized as follows. In section 2, we present our model of shared-memory systems, and provide a formal definition of conditional primitives. In Section 3 we prove space lower bounds for long-lived objects. In Section 4 we prove time/space tradeoffs for implementations of one-time objects. Section 5 concludes with a discussion of the results.

## 2. PRELIMINARIES

### 2.1 Shared Memory System Model

Our model of an asynchronous shared memory system is largely based on the model described by Cypher in [4], which is based, in turn, on the model given by Merritt and Taubenfeld [11]. Shared objects are specified as in [7]. An object is an abstract data type. It is characterized by a set of possible values and by a set of *operations* that provide the only means to manipulate it. A shared memory protocol provides a specific data-representation for the object, and a specific implementation for each of its operations. An  $n$ -process shared memory protocol consists of a non-empty set  $\mathbf{E}$  of executions, a set  $\mathbf{P}$  of  $n$  processes, a set  $\mathbf{R}$  of shared memory registers, and a function  $\mathbf{I}$  that assigns an initial value to each register in  $\mathbf{R}$ . Any register may support one or more operations including Read, Write, and various types

of atomic Read-Modify-Write operations. No bound is assumed on register size (i.e. the number of different possible values the register can have). An *execution-fragment* is a sequence (either finite or infinite) of *events*, where an event is an operation that a specific process applies to a specific register. An *execution* is an execution-fragment that starts from the initial state. An event  $e$  can have one of the following three forms:

- $\text{Read}(p,r)$ : indicates that process  $p$  atomically reads the value of register  $r$ ;
- $\text{Write}(p,r,w)$ : indicates that process  $p$  atomically writes the value  $w$  to register  $r$ ;
- $\text{RMW}(p,r,w,g,h)$ : indicates that the code shown in Figure 1 (a) is atomically performed on behalf of process  $p$ . The write function,  $g$ , and the return-value function,  $h$ , both receive two parameters,  $v$ , the value of register  $r$  when event  $e$  starts to execute, and  $w$ , the input value to the RMW event. The type of a RMW operation is specified by the ordered pair  $\langle g,h \rangle$  consisting of its write and return-value functions.<sup>1</sup>

As an example, the write and return-value functions of the CAS operation are shown in Figure 1 (b).

<pre> RMW(p,r,w,g,h) v=r.Read(); r.Write(g(v,w)); return h(v,w); </pre> <p>(a)</p>	<pre> CAS.g(v,w=&lt;old,new&gt;) if (v==w.old)     return w.new; else     return v;  CAS.h(v,w=&lt;old,new&gt;) if (v==w.old)     return true; else     return false; </pre> <p>(b)</p>
--	---

**Figure 1: (a) The atomic operation of a Read-Modify-Write event. (b) The write and return-value functions of the CAS operation. The input of the CAS operation,  $w$ , is a compound value.**

We next define the class of conditional RMW primitive operations we investigate in this paper.

**DEFINITION 2.1.** *Let  $Op = \langle g,h \rangle$  be a RMW operation. We say that  $Op$  is a conditional operation if, for every possible input value  $w$ , the following holds:*

$$\left| \{v \mid g(v,w) \neq v\} \right| = 1.$$

*Let  $Op = \langle g,h \rangle$  be a conditional operation and let  $e = \text{RMW}(p,r,w,g,h)$  be an event of type  $Op$ . We call the single value  $v_w \neq g(v_w,w)$  the change-point of  $e$ . We call any other value  $v$ ,  $v \neq v_w$ , a fixed-point of  $e$ .*

In other words, a RMW operation is conditional if, for every input  $w$ , there is exactly one value  $v_w$  such that  $g(v_w,w) \neq v_w$ . If  $g$  and  $w$  are such that  $g(v,w) = v$  for all values of  $v$ , then the event  $\text{RMW}(p,r,w,g,h)$  is essentially an event that never writes. For any event  $e = \text{RMW}(p,r,w,g,h)$ ,

<sup>1</sup>Our definition of a RMW operation-type is a generalization of the definition of comparison primitives that appears in Section 5 of [1].

the value  $v_w$  is the only value of  $r$  that will be changed by  $e$ . We observe that CAS and Test-and-set are conditional operations by Definition 2.1. Following [9], we model the Load-Linked (called *LL*) and Store-conditional (called *SC*) operations as follows. The state of every shared register is a structure,  $r$ , that contains two fields:  $r.value$  and  $r.processSet$ . When some process,  $p$ , performs an *LL* operation on some register  $r$ ,  $p$  is added to  $r.processSet$  and *LL* returns  $r.value$ . When  $p$  performs an *SC*( $v$ ) operation on  $r$ , then there are the following two possibilities: (1) If  $p \in r.processSet$ , the *SC* operation sets  $r.processSet \leftarrow \phi$ , sets  $r.value \leftarrow v$  and returns  $(true, v)$ . In this case we say the *SC* was successful. (2) Otherwise, *SC* returns  $(false, r.value)$ , in which case we say the *SC* was unsuccessful. Modelled this way, it is easily seen that *LL* is equivalent to a Read and *SC* is a conditional by Definition 2.1.

We call a register that supports only Read and Write operations a *read-write register*; we call a register that supports only conditional operations a *conditional register*; we call a register that supports Read, Write, and one or more conditional operations a *read-write-conditional register*.

We denote the empty execution by  $\epsilon$ . Given any event  $e$ , we say that  $e$  is a *writing event* if  $e$  is a Write event or if  $e$  is a RMW event; if  $e$  is a Read event or if  $e$  is a RMW event, we say that  $read(e)$  holds.

The memory register *accessed* (read and/or written) by  $e$  is denoted  $mem(e)$ . The process that executes  $e$  is denoted  $proc(e)$ . For any two execution-fragments  $E$  and  $E'$ , where  $E$  is finite, the execution-fragment  $E \circ E'$  denotes the concatenation of  $E$  and  $E'$ . For any execution-fragment  $E$ , we define  $procs(E)$  to be the set of processes that perform some event in  $E$ . Let  $r \in R$  be a memory register and let  $E \in \mathbf{E}$  a finite execution, then  $value(E,r)$  denotes the last value written by an event in  $E$  to  $r$ , or  $\mathbf{I}(r)$  if there was no such event. Given an execution  $E$  and any subset  $P \subseteq \mathbf{P}$ , we let  $proj(E,P)$  denote only those events in  $E$  that were performed by processes in  $P$ . If  $P = \{p\}$ , we also use the notation  $proj(E,p)$  instead of  $proj(E,\{p\})$ . Let  $E \in \mathbf{E}$  be a finite execution, and  $e$  be an event. If  $E \circ e \in \mathbf{E}$  holds, we say that  $e$  is *enabled after*  $E$ . An execution-fragment  $E$  is *fair*, if any process that has an enabled event just before  $E$  starts is in  $procs(E)$ . An execution is fair if all of its suffixes are fair.

## 2.2 High Level Operations

The execution of a high-level object operation involves, in general, both private- and shared-memory events; in this paper we only deal with shared-memory events. Thus, we view an execution of a high-level operation  $Op$  as consisting of a sequence of one or more atomic shared memory events, each of which can be either *Read*, *Write*, or *RMW*. Let  $Op$  be an execution of a high-level operation performed by some process in some execution  $E$ . We assume processes execute at most a single high-level operation at any given time. We denote by  $proc(E,Op)$  the process that executes  $Op$  in  $E$ ; we denote by  $events(E,Op)$  the sequence of memory-events performed by  $proc(Op)$  while executing  $Op$  in  $E$ . Whenever  $E$  is clear from the context, we simply write  $proc(Op)$  and  $events(Op)$ . We say that a process  $p$  is *active* after  $E$ , and write  $active(E,p)$ , if, after  $E$ ,  $p$  is in the middle of executing some high-level operation  $Op$ , i.e.  $p$  performed at least one event of an instance of  $Op$ , but has not performed the last event of this instance of  $Op$ . If  $p$  is not active after  $E$ , we say

that  $p$  is *idle* after  $E$ , and write  $idle(E, p)$ . We say that an execution  $E$  is *quiescent* if all processes are idle after  $E$ . We say that an operation  $Op$  is contained within an execution  $E$ , and write  $Op \subset E$ , if  $Op$  terminates in  $E$ . If  $Op \subset E$ , we denote by  $result(Op, E)$  the value returned by  $Op$  in  $E$ . We say that a high-level operation  $Op$  is enabled after an execution  $E$ , and write  $enabled(E, Op)$ , if  $proc(Op)$  has an enabled event  $e$  of  $Op$  after  $E$ . Assume  $enabled(E, Op)$  holds. Then we denote by  $read\_solo(E, Op)$  the (possibly empty) execution fragment, consisting of Read events, that results when we let  $proc(Op)$  run solo after  $E$ , until either it is about to perform a writing event, or it returns from  $Op$ .

For presentation simplicity, we often refer to an underlying execution as a *state*. Thus, for example, instead of saying that a high-level operation  $Op$  is enabled after execution  $E$ , sometimes we say that after  $E$  the system is in a state where  $Op$  is enabled. If  $S$  is the system state after  $E$ , we also use the notation  $read\_solo(S, Op)$  instead of  $read\_solo(E, Op)$ .

### 3. LONG-LIVED OBJECTS

In this section, we obtain lower bounds on the memory requirements of wait-free implementations of long-lived objects in a class of objects that we call *Visible(n)*, when implementations use registers that only support Read, Write and/or conditional operations. We also consider starvation-free implementations of mutual exclusion from such registers.

First, we define the key concepts of *invisible* and *visible* events. Informally, an invisible event is a writing event by some process that cannot be observed by other processes. An invisible event was termed an *obliterated* event in [3]. It was defined there only for read-write registers. We extend this notion to conditional RMW events.

The following two definitions formally define the concepts of invisible RMW and Write events.

**DEFINITION 3.1.** *Let  $e$  be a RMW( $p, r, w, g, h$ ) event in an execution  $E$ , where  $E = E_1 \circ e \circ E_2$ . We say that  $e$  is invisible in  $E$ , and write  $invisible(E, e)$ , if at least one of the following holds:*

- $g(value(E_1, r), w) = value(E_1, r)$ ;
- $E_2 = e' \circ E_3$ ,  $e'$  is a Write event and  $mem(e') = r$ .

In other words, a RMW event  $e$  is invisible, if either the value of the register on which  $e$  operates is a fixed-point of  $e$ , or if  $e$  is immediately followed by a Write event on the same register. We note that  $e$  would also be made invisible to other processes in an execution where it is followed by a sequence of events on other registers that terminates with a Write event applied to  $r$ ; however our proofs do not require this broader definition.

**DEFINITION 3.2.** *Let  $e = Write(p, r, w)$  be an event in an execution  $E$  where  $E = E_1 \circ e \circ E_2 \circ e' \circ E_3$ . We say that  $e$  is invisible in  $E$ , and write  $invisible(E, e)$ , if the following holds:*

- $e'$  is a Write event and  $mem(e') = r$ ;
- $\forall f \in E_2 : \neg(read(f) \wedge (mem(f) = r) \wedge (proc(f)) \neq proc(e))$ .

In other words, a Write event  $e$ , issued by some process  $p$ , is invisible in an execution  $E$ , if it is overwritten in  $E$  by another Write event before it is read by another process. Note that whereas Write events can only be made invisible by *subsequent* Write events, RMW events can be made invisible either by *preceding* writing events, or by an immediately following Write event. If  $e$  is a Write or RMW event that is not invisible in  $E$ , we say that  $e$  is a *visible* event in  $E$ .

We next define the class *Visible(n)*, to which our results apply.

**DEFINITION 3.3.** *An object  $O$  is in  $Visible(n)$ , if  $O$  supports a high-level operation  $Op$  such that for every wait-free protocol,  $P$ , implementing  $O$  and every execution  $E \in \mathbf{E}$  of  $P$ ,*

- $Op$  is always enabled at all  $n$  processes, and
- For every execution  $OP_i$  of  $Op$  such that  $OP_i \subset E$ ,  $events(E, OP_i)$  includes at least one visible writing event.

It is easily proven that many key concurrent objects are in *Visible(n)*: counters, stacks, queues, swap, fetch-and-add, and single-writer snapshot are some examples. As one example, we prove in the following that concurrent counting is in *Visible(n)*.

A concurrent counter object supports a single fetch-and-increment (*FAI*) operation. The counter values returned by different *FAI* operations in an execution are required to be distinct natural numbers. It is also required that, in any quiescent state, the values returned by the counter constitute a contiguous range of natural numbers.<sup>2</sup>

**LEMMA 3.1.** *Concurrent counting is in  $Visible(n)$ .*

**PROOF.** Assume otherwise by way of contradiction. Let  $P$  be a wait-free protocol implementing a concurrent counter. By assumption, there is at least one execution  $E$  of  $P$  such that some *FAI* operation  $F \subset E$ , performed by some process  $p$ , terminates in  $E$  and does not perform a visible write. We assume without loss of generality that  $E$  is quiescent. (Otherwise, as  $P$  is wait-free, we can extend  $E$  by letting all active processes run until they finish their high-level operations.) Let  $V$  be the number of *FAI* operations completed in  $E$  (including  $F$ ). Then the counter must have returned the values  $1, \dots, V$ . We now extend  $E$  by a solo execution-fragment,  $E_1$ , of another *FAI* operation,  $F_1$ , by process  $q \neq p$ .  $F_1$  must return value  $V + 1$ . However, by assumption, there is another execution of  $P$ ,  $E'$ , such that  $p$  has no events in  $E'$  and the following holds:

$$\forall p' \neq p : proj(E', p') = proj(E, p').$$

Consequently, there exists an execution  $E_2 = E' \circ F_1$  and  $F_1$  returns  $V + 1$  also in  $E_2$ . Note that in  $E'$  the counter has returned only  $V - 1$  values; it follows that, after  $E_2$  terminates, we are in a quiescent state and the values returned by the counter do not constitute a contiguous range, contradicting the definition of a counter.  $\square$

<sup>2</sup>*Linearizable concurrent counters* are also required to be linearizable [8], i.e. if  $FAI_i$  and  $FAI_j$  are two executions of the *FAI* operation, and  $FAI_i$  entirely precedes  $FAI_j$  in an execution  $E$ , then  $result(FAI_i, E) < result(FAI_j, E)$  must hold. The proof of Lemma 3.1 applies for both linearizable and non-linearizable concurrent counters, thus both counter types are in *Visible(n)*.

The following definition is key to our proofs.

**DEFINITION 3.4.** Let  $A = e_1, e_2, \dots, e_k$  be a sequence of  $k$  writing events by  $k$  different processes, pending in some state  $S$ . We say that  $S$  is a  $k$ -levelled state with  $e_j$  at level  $j$  for  $j = 1, \dots, k$ , if, in  $S$ , events of higher levels cannot make events of lower levels invisible. More formally, for all levels  $j = 1, \dots, k$  and for every ordering  $\sigma_I$  of every subsequence  $I$  of  $A$ ,  $I = e_{i_1}, \dots, e_{i_{|I|}}$ , the following holds:

$$(e_j \in I) \wedge (\forall i < j : e_i \notin I) \implies \text{visible}(S \circ e_{\sigma(i_1)} \circ \dots \circ e_{\sigma(i_{|I|})}, e_j)$$

We call  $A$  a  $k$ -levelled-sequence in  $S$ .

The organization of the proofs for long-lived objects is as follows. We first prove that all wait-free implementations of objects in  $Visible(n)$  and all starvation-free implementations of mutual exclusion can be brought to an  $n$ -levelled state. We then prove that any protocol that uses only read-write registers or only conditional registers, and reaches a  $k$ -levelled state, uses at least  $k$  such registers, and that any protocol that uses only read-write-conditionals registers, and reaches a  $k$ -levelled state, uses at least  $\lceil \frac{k}{2} \rceil$  such registers. Combining these results, we obtain the space lower bounds for long-lived objects.

The following lemma proves that wait-free implementations of objects in  $Visible(n)$  can be made to reach  $n$ -levelled states.

**LEMMA 3.2.** Let  $P$  be a wait-free protocol implementing some object  $O \in Visible(n)$ . Then  $P$  can be brought to an  $n$ -levelled state.

**PROOF.** In the following, we denote the execution of a high-level operation  $Op$  by process  $p_i$ ,  $1 \leq i \leq n$ , as  $Op_i$ . We construct an execution  $E$  incrementally, in  $n$  phases. In phase  $i$ , we extend the execution so that  $p_i$  is on the verge of performing the first visible writing event (within its current high-level operation) that cannot be made invisible by processes  $p_{i+1}, \dots, p_n$ . Figure 2 shows pseudo-code describing the construction of  $E$ . We next prove that the following claims hold for each phase  $i$ ,  $1 \leq i \leq n$ .

1. The construction of phase  $i$  terminates;
2. Right after the termination of phase  $i$ , each of the processes  $p_1, \dots, p_i$  has a pending writing event,  $e_i$ , that cannot be made invisible by events of processes  $p_{i+1}, \dots, p_n$ ; moreover, none of the operations  $Op_1, \dots, Op_i$  has performed a visible writing event yet;
3. Right after the termination of phase  $i$ , processes  $p_{i+1}, \dots, p_n$  are idle.

The proof proceeds by induction on the phase number,  $i$ . Clearly, before the first phase all of the above claims hold vacuously. We now assume the claims hold after the termination of phase  $i - 1$ , and prove for phase  $i$ . We first show that the construction of phase  $i$  terminates. From wait-freedom, the read-solo execution by process  $p_i$  (line 4) is finite. By induction hypothesis, before phase  $i$  begins,  $p_i$  is idle, and so, as  $O$  is in  $Visible(n)$ ,  $Op_i$  must perform a visible write in the course of phase  $i$ . Consequently, the read-solo execution by  $p_i$  must terminate with a pending writing event, which we denote  $e_i$ . Next, we show that the phase construction eventually reaches line 9. There are two cases to consider:

```

1:  $E \leftarrow \epsilon$ 
2: for ( $i = 1$ ;  $i \leq n$ ;  $i++$ ) // Construct  $E$  in  $n$  phases.
3: {
  construct-read-solo:

4:  $E \leftarrow E \circ \text{read\_solo}(E, Op_i)$ 
5: Let  $e_i$  denote the pending writing event by  $p_i$ 
6: if exist  $E_1, E_2$  by  $\{p_{i+1}, \dots, p_n\}$  s.t.
    $\text{invisible}(E \circ E_1 \circ e_i \circ E_2, e_i)$  holds
   {
7:    $E \leftarrow E \circ E_1 \circ e_i \circ E_2$ 
8:   goto construct-read-solo
   }
9: else if ( $i < n$ )
10:  if any of  $Op_{i+1}, \dots, Op_n$  are enabled
11:    $E \leftarrow E \circ E_3$ , where in  $E_3$   $Op_{i+1}, \dots, Op_n$ 
    run to completion
12: }
```

**Figure 2: Pseudo-code for the constructing of execution  $E$ .**

- If there is no extension  $E_1 \circ e_i \circ E_2$  of  $E$ , with  $\text{procs}(E_1), \text{procs}(E_2) \subseteq \{p_{i+1}, \dots, p_n\}$  that makes  $e_i$  invisible, then the construction reaches line 9 immediately.
- Otherwise, there is an extension  $E_1 \circ e_i \circ E_2$  (possibly with an empty  $E_1$  and/or an empty  $E_2$ ), with  $\text{procs}(E_1), \text{procs}(E_2) \subseteq \{p_{i+1}, \dots, p_n\}$ , that makes  $e_i$  invisible. We set  $E \leftarrow E \circ E_1 \circ e_i \circ E_2$  (line 7) and jump to *construct-read-solo*. As  $P$  is wait-free, and as every jump to *construct-read-solo* implies that  $Op_i$  performs at least one additional event, we can only jump to *construct-read-solo* a finite number of times, after which we must reach line 9.

Having reached line 9, if  $i < n$ , we extend  $E$  with the execution-fragment  $E_3$ , where in  $E_3$  all of the active processes in  $\{p_{i+1}, \dots, p_n\}$ , if any, complete their high-level operations. The existence of such a finite execution-fragment  $E_3$  is guaranteed from the wait-freedom of  $P$ ; thus, claims 1 and 3 are proven for phase  $i$ . Claim 2 follows directly from the construction and from the induction hypothesis. Let  $S$  be the state after the execution of  $E$ . By construction, for every  $1 \leq i \leq n$ ,  $\text{visible}(S \circ E_1 \circ e_i \circ E_2, e_i)$  holds for any execution-fragments  $E_1, E_2$  consisting of any ordering of any disjoint subsets of the events  $e_{i+1}, \dots, e_n$ . Thus  $A = e_1, \dots, e_n$  is an  $n$ -levelled-sequence in  $S$ , with  $e_i$  having level  $i$ , and  $S$  is an  $n$ -levelled state.  $\square$

Next we show the same result for starvation-free mutual exclusion. We model mutual exclusion as follows. A mutual exclusion protocol supports a single operation, called *Mutex*, which has the structure shown in Figure 3. The non-critical section is assumed to take place outside of *Mutex* operations. It is assumed that no process halts while executing the *Mutex* operation. The requirements from a mutual exclusion implementation are the following:

- *Exclusion*: at most one process executes its critical section at any time.

```

Mutex()
{
  Entry Section
  Critical Section
  Exit Section
}

```

**Figure 3: Structure of the mutual-exclusion operation.**

- *Starvation-freedom* (also known in the literature as *lockout-freedom*): in any fair execution, if some process is in its entry section, then that process eventually executes its critical section.
- *Finite exit*: in any fair execution, if a process is in its exit section, then it eventually exits the Mutex operation.

We need the following technical lemma.

LEMMA 3.3. *Let  $P$  be a starvation-free mutual exclusion protocol, and let  $E$  be an execution of  $P$  where some process  $p$  enters the critical section, then  $E$  contains a visible writing event by  $p$ .*

PROOF. Assume otherwise by way of contradiction, then there is some execution  $E$  after which  $p$  is in the critical section, and  $E$  does not contain any visible write by  $p$ . Consequently, there exists another execution,  $E'$ , that does not contain any events by process  $p$ , such that:

$$\forall p' \neq p : \text{proj}(E', p') = \text{proj}(E, p').$$

From starvation-freedom, there is some extension  $E''$  of  $E'$ , such that some process  $q \neq p$  enters the critical section in  $E' \circ E''$  and  $p \notin \text{procs}(E' \circ E'')$ . Thus,  $q$  enters the critical section also in the execution  $E \circ E''$ , and exclusion is violated, a contradiction.  $\square$

LEMMA 3.4. *Let  $P$  be a long-lived starvation-free mutual exclusion protocol. Then  $P$  can be brought to an  $n$ -levelled state.*

The proof of Lemma 3.4 is almost identical to that of Lemma 3.2. An execution that results in an  $n$ -levelled state is constructed in exactly the same way. The only differences between the two proofs are that here we use starvation-freedom instead of wait-freedom and we use Lemma 3.3 instead of membership in  $\text{Visible}(n)$ .

We next prove that any protocol that reaches a  $k$ -levelled state and uses registers that support only Read, Write and/or conditional operations, uses  $\Omega(k)$  such registers.

LEMMA 3.5. *Assume that a protocol  $P$  can be brought to a  $k$ -levelled state  $S$  for some  $k > 0$ . Let  $\text{SPACE}(P)$  denote the number of registers used by  $P$ .*

1. *If  $P$  uses only read-write registers, then  $\text{SPACE}(P) \geq k$ .*
2. *If  $P$  uses only conditional registers, then  $\text{SPACE}(P) \geq k$ .*

3. *If  $P$  uses only read-write-conditional registers, then  $\text{SPACE}(P) \geq \lceil \frac{k}{2} \rceil$ .*

PROOF. Let  $S$  be a  $k$ -levelled state with a  $k$ -levelled sequence  $A = e_1 \cdots e_k$ . Let  $p_i = \text{proc}(e_i)$  and let  $R = \{\text{mem}(e_i) | i = 1, \dots, k\}$ . Note that no two Write events can be pending on the same register because any of these two events makes the other invisible. Thus, a lower-level event can be made invisible by a higher-level event, contradicting our assumption that  $A$  is a  $k$ -levelled sequence of  $S$ . This proves (1).

Next we show that no two conditional RMW events can be pending on the same register. Assume otherwise. Let  $e_i = \text{RMW}(p_i, r, w_i, g_i, h_i)$  and  $e_j = \text{RMW}(p_j, r, w_j, g_j, h_j)$  be two conditional RMW events pending on the same register  $r \in R$ . Since  $A$  is a  $k$ -levelled sequence in  $S$ , both  $\text{visible}(S \circ e_i, e_i)$  and  $\text{visible}(S \circ e_j, e_j)$  hold. Consequently  $\text{value}(S, r)$  is the only change-point of both  $e_i$  and  $e_j$ . This implies that either of the events  $e_i$  or  $e_j$  can make the other invisible, contradicting our assumption that  $A$  is a  $k$ -levelled sequence of  $S$ . This proves (2).

No two Write events can be pending on the same register and no two conditional RMW events can be pending on the same register. It follows that at most two events can be pending on any single register  $r \in R$ : one Write and one conditional RMW. This proves (3).  $\square$

Our main result follows directly from the above lemmata.

THEOREM 3.6. *Let  $P$  be an  $n$ -process protocol that is a starvation-free implementation of mutual exclusion or a wait-free implementation of an object in  $\text{Visible}(n)$ . Let  $\text{SPACE}(P)$  denote the number of registers used by  $P$ .*

- *If  $P$  uses only read-write registers, then  $\text{SPACE}(P) \geq n$ .*
- *If  $P$  uses only conditional registers, then  $\text{SPACE}(P) \geq n$ .*
- *If  $P$  uses only read-write-conditional registers, then  $\text{SPACE}(P) \geq \lceil \frac{n}{2} \rceil$ .*

PROOF. Immediate from Lemmata 3.2, 3.4 and 3.5  $\square$

## 4. ONE-TIME OBJECTS

A one-time object is an object whose operations can be called by a process only once. It is easily seen that the memory bounds proven in Section 3 for long-lived objects do not hold, in general, for one-time objects. As an example, a one-time wait-free counter shared by  $n$  processes can be implemented by using a single register that supports the *CAS* operation. However, by Theorem 3.6, a long-lived wait-free counter requires at least  $n$  such registers. For one-time objects, we prove time-space tradeoffs for wait-free implementations that use only read-write-conditional registers. Our proofs apply for a class of objects defined as follows.

DEFINITION 4.1. *An object  $O$  is in  $\text{Visible}'(n)$ , if  $O$  supports a high-level operation  $Op$  such that for every wait-free protocol,  $P$ , implementing  $O$ ,*

- *in the initial state, all  $n$  processes have an enabled  $Op$  operation and*

- for every execution  $E \in \mathcal{E}$  of  $P$  and every execution-instance  $Op_i$  of  $Op$ ,  $Op_i \subset E$ ,  $events(E, Op_i)$  includes at least one visible writing event.

It is easily seen that one-time counters, stacks, queues, swap, fetch-and-add, and single-writer snapshot objects are in  $Visible'(n)$ . We note that the following tradeoffs also apply to starvation-free implementations of one-time mutual exclusion.

We first need the following lemma, that states that when multiple conditional events and Write events are pending on the same register, they can always be scheduled such that at most two events are visible.<sup>3</sup>

LEMMA 4.1. *Let  $L = \{e_1, \dots, e_k\}$  be a set of  $k$  writing events pending on the same register  $r$  in state  $S$ ; then, if all the events are either Write or conditional events, they can be scheduled such that at most two of them are visible.*

PROOF. If all the events in  $L$  are Write events, the lemma is immediate, as only the Write event scheduled last may be visible. Otherwise, there are some conditional events in  $L$ . Let  $L_1$  contain all Write events in  $L$ , and let  $L' = L \setminus L_1$  be the non-empty set of conditional events in  $L$ . Let  $cur$  be the value of  $r$  in state  $S$ , and consider any event  $e \in L'$ ,  $e = RMW(p, r, w, g, h)$ . As  $\langle g, h \rangle$  is conditional, there is only a single value  $v$  such that  $g(v, w) \neq v$ ; we denote this value by  $old(e)$ . We partition the events of  $L'$  to the following two disjoint sets:  $L' = L_2 \cup L_3$ , where  $L_2 = \{e_j \in L' \mid old(e_j) \neq cur\}$ , and  $L_3 = L' \setminus L_2$ . In any schedule where the events of  $L_2$  are performed first (in any order) followed by all the events of  $L_3$  (in any order), we have at most a single visible event (namely the first event scheduled from  $L_3$ , if any). If  $L_1$  is non-empty, we schedule the events in it last (in any order), and the last of these to be scheduled may also be visible.  $\square$

The following theorem states a tradeoff between the number of read-write-conditional registers used by a protocol and the number of write events it may have to perform.

THEOREM 4.2. *Let  $O$  be an object in  $Visible'(n)$ , and let  $P$  be an  $n$ -process wait-free implementation of  $O$  that uses only read-write-conditional registers. If  $P$  uses at most  $m < n$  registers, then  $P$  has an execution  $E$  in which the total number of writing events performed is in  $\Omega(\frac{n^2}{m})$ .*

PROOF. For simplicity, we assume that  $2 \cdot m$  divides  $n$ . As  $O$  is in  $Visible'(n)$ , the high-level operations performed by the  $n$  processes cannot terminate before they perform a visible write. We construct an execution  $E$  in phases. In each phase, we bring all processes to be on the verge of writing. As  $P$  can only apply Writes and conditional RMW operations to registers, from Lemma 4.1, an adversarial scheduler can arrange the writing events on each register  $r$  so that at most two of these events are visible in  $E$  at  $r$ . Consequently, at most  $2m$  processes perform a visible write during the phase. We let these processes run to completion. In the second phase there remain  $n - 2m$  processes that did not perform a visible write. The same argument can be applied

<sup>3</sup>A claim very similar to Lemma 4.1 appears in Section 5 of [1] without proof.

to them. We continue in this manner, until all processes have performed a visible write. Denoting the total number of writes made in  $E$  by  $W(P)$ , we get:

$$W(P) = \sum_{k=0}^{\frac{n}{2m}} (n - 2km) \in \Omega(\frac{n^2}{m}). \quad (1)$$

$\square$

From Theorem 4.2 we get the following, more specific, tradeoff between space and number of write operations:

THEOREM 4.3. *Let  $P$  be an  $n$ -process wait-free implementation of  $O$  that uses only read-write-conditional registers and let  $SPACE(P)$ ,  $W_{amortized}(P)$  and  $W_{op}(P)$ , respectively, denote the number of registers used by  $P$ , the worst-case amortized number of writing events issued by  $n$  concurrent operations, and the worst-case number of writing events issued by a single high-level operation in  $P$ . Then either  $SPACE(P) > \sqrt{n}$  or both  $W_{amortized}(P)$  and  $W_{op}(P)$  are in  $\Omega(\sqrt{n})$ .*

The results of Theorems 4.2 and 4.3 can be strengthened, by counting not only the number of writes, but also the number of *memory stalls* caused from write contention. In all shared-memory systems, when multiple processes attempt to write to the same memory location simultaneously, the writes are serialized, and waiting operations incur memory stalls. The concept of memory stalls was introduced in [5]. We use the following definition of memory stalls, which is stricter than that of [5], as it counts only stalls caused by contention in *writing*.

DEFINITION 4.2. *Let  $E$  be an execution and let  $e_0, \dots, e_l$  be a maximal sequence of two or more consecutive events in  $E$  such that*

1.  $e_j$  is either a Write or a RMW event, for  $j \in \{0, \dots, l\}$ , and
2.  $proc(e_{j_1}) \neq proc(e_{j_2})$  and  $mem(e_{j_1}) = mem(e_{j_2})$  for distinct  $j_1, j_2$  in  $\{0, \dots, l\}$ .

*Let  $Op$  be the high-level operation whose execution issued  $e_j$ . Then we say that  $Op$  incurs  $j$  memory stalls in  $E$  on account of  $e_j$ .*

We obtain a tradeoff between space complexity and the number of memory stalls incurred by a wait-free implementation of any object in  $Visible'(n)$  that uses only read-write-conditional registers.

THEOREM 4.4. *Let  $P$  be an  $n$ -process wait-free implementation of  $O$  that uses only read-write-conditional registers. If  $P$  uses at most  $m < n$  registers, then there is an execution in which the total number of memory stalls incurred by all processes is in  $\Omega(\frac{n^3}{m^2})$ .*

PROOF. For simplicity, and without loss of generality, we assume that  $2 \cdot m$  divides  $n$ . We consider the same type of phased execution  $E$  we constructed in Theorem 4.2. We denote by  $MS(E)$  the total number of memory stalls incurred by processes during the execution of  $E$ ; we denote by  $M_k(E)$ ,

for  $1 \leq k \leq n/2m$ , the total number of memory stalls incurred by processes on account of their writing events in phase  $k$ , and for  $1 \leq i \leq m$  we denote by  $Writers(E, k, i)$  the number of processes that write to register  $r_i$  in phase  $k$  of  $E$ . We get:

$$M_k(E) = \sum_{i=1}^m \frac{Writers(E, k, i) \cdot (Writers(E, k, i) - 1)}{2} \quad (2)$$

$$\sum_{i=1}^m Writers(E, k, i) = n - 2 \cdot k \cdot m. \quad (3)$$

It is easily seen that  $M_k(E)$  is minimized in the equations above when the writes are equally distributed among registers (i.e. every register is written by  $\frac{n-2 \cdot k \cdot m}{m}$  processes), thus we get:

$$M_k(E) \geq \frac{m}{2} \cdot \left(\frac{n-2km}{m}\right) \cdot \left(\frac{n-2km}{m} - 1\right). \quad (4)$$

Summing over all phases, we get:

$$\begin{aligned} MS(E) &= \sum_{k=0}^{\frac{n}{2m}} M_k(E) \\ &\geq \frac{m}{2} \cdot \sum_{k=0}^{\frac{n}{2m}} \left(\frac{n-2km}{m}\right) \cdot \left(\frac{n-2km}{m} - 1\right) \in \Omega\left(\frac{n^3}{m^2}\right). \end{aligned} \quad (5)$$

□

From Theorem 4.4 we get the following, more specific, tradeoff between space and number of memory steps:

**THEOREM 4.5.** *Let  $P$  be an  $n$ -process wait-free implementation of  $O$  that uses only read-write-conditional registers. Let  $SPACE(P)$ ,  $MS_{amortized}(P)$  and  $MS_{op}(P)$ , respectively, denote the maximal number of registers used by  $P$ , the worst-case amortized number of memory stalls incurred by  $n$  concurrent operations, and the worst-case number of memory stalls incurred by a single high-level operation in  $P$ . Then either  $SPACE(P) > n^{\frac{2}{3}}$  or both  $MS_{amortized}(P)$  and  $MS_{op}(P)$  are in  $\Omega(n^{\frac{2}{3}})$ .*

## 5. DISCUSSION

Conditional synchronization primitives are among the strongest primitives according to Herlihy's wait-free hierarchy, as they can be used to implement deterministic wait-free consensus for any number of processes. This paper shows that conditional synchronization primitives are relatively inefficient in terms of memory space for wait-free implementations of most non-trivial objects. Our results apply to starvation-free implementations of mutual exclusion and to wait-free implementations of objects in  $Visible(n)$ , a class that contains objects supporting an operation that must perform a visible write before it terminates. In contrast, starvation-free mutual exclusion and some of the key objects in  $Visible(n)$  can be implemented using  $O(1)$  registers that support other synchronization primitives, such as fetch-and-increment.

Several researchers have previously obtained results that indicate the weakness of conditional primitives. Anderson and Kim [1], continuing the work of Cypher [4], proved a  $\Omega(\log n / \log \log n)$  remote memory references lower bound on starvation-free mutual exclusion implementations that use CAS and LL/SC, "almost" equal to the  $\log n$  remote

memory references upper bound that uses read-write registers. In contrast, starvation-free mutual exclusion can be implemented in  $O(1)$  remote memory references by using registers that support fetch-and-increment. Jayanti [9] considers wait-free object implementations that use registers that support read/write and any of the following operations: LL/SC, *swap*, and *move*. He proves a lower bound of  $\Omega(\log n)$  on the worst-case latency of any such implementation for counters, stacks and queues. This is yet another indication of the weakness of LL/SC in the context of wait-free implementations.

*Multi-valued* objects are objects such that different operations (called by different processes) may return different, yet related, values. Consensus is the key concurrent object that is not multi-valued. The aforementioned results, viewed collectively, imply that even though CAS and LL/SC are strong in the context of solving deterministic wait-free consensus, they are relatively weak in the context of wait-free implementations of most multi-valued objects. The conclusion is therefore, that basing multiprocessor synchronization solely on conditional synchronization primitives is probably not a good design choice.

## 6. REFERENCES

- [1] J. Anderson and Y. Kim. An improved lower bound for the time complexity of mutual exclusion. In *ACM Symposium on Principles of Distributed Computing*, pages 90–99, 2001.
- [2] J. Anderson and Y. Kim. Shared-memory mutual exclusion: Major research trends since 1986. *Distributed Computing*, 16(2-3):75–110, 2003.
- [3] J. Burns and N. Lynch. Bounds on shared memory for mutual exclusion. *Information and Computation*, 107(2):171–184, December 1993.
- [4] R. Cypher. The communication requirements of mutual exclusion. In *ACM Proceedings of the Seventh Annual Symposium on Parallel Algorithms and Architectures*, pages 147–156, 1995.
- [5] C. Dwork, M. Herlihy, and O. Waarts. Contention in shared memory algorithms. *Journal of the ACM (JACM)*, 44(6):779–805, 1997.
- [6] F. E. Fich and E. Ruppert. Hundreds of impossibility results for distributed computing. *Distributed Computing*, 16:121–163, 2003.
- [7] M. Herlihy. Wait-free synchronization. *ACM Transactions On Programming Languages and Systems*, 13(1):123–149, Jan. 1991.
- [8] M. P. Herlihy and J. M. Wing. Linearizability: A correctness condition for concurrent objects. *ACM Transactions On Programming Languages and Systems*, 12(3):463–492, July 1990.
- [9] P. Jayanti. A time complexity lower bound for randomized implementations of some shared objects. In *Symposium on Principles of Distributed Computing*, pages 201–210, 1998.
- [10] P. Jayanti, K. Tan, and S. Toueg. Time and space lower bounds for non-blocking implementations. In *Symposium on Principles of Distributed Computing*, pages 257–266, 1996.
- [11] M. Merrit and G. Taubenfeld. Knowledge in shared memory systems. In *ACM Symp. on Principles of Distributed Computing*, pages 189–200, 1991.