

Computing with Infinitely Many Processes under assumptions on concurrency and participation*

Michael Merritt[†]

Gadi Taubenfeld[‡]

February 19, 2006

Abstract

We explore four classic problems in concurrent computing (election, mutual exclusion, consensus, and naming) when the number of processes which may participate is infinite. Partial information about the number of actually participating processes and the concurrency level is shown to affect the possibility and complexity of solving these problems. We survey and generalize work carried out in models with finite bounds on the number of processes, and prove several new results. These include improved bounds for election when participation is required (even for finitely many processes, as investigated by Styer and Peterson [SP89]) and a new adaptive starvation-free mutual exclusion algorithm for unbounded concurrency. We survey results in models with shared objects stronger than atomic registers, such as test&set bits, semaphores or read-modify-write registers, and update them for the infinite process case.

*A preliminary version of the results presented in this paper appeared in *Proceedings of the 12th International Symposium on Distributed Computing*, LNCS 1914:164-178, October 2000.

[†]AT&T Labs, 180 Park Ave., Florham Park, NJ 07932-0971. mischu@research.att.com.

[‡]The Interdisciplinary Center, P.O.Box 167, Herzliya 46150, Israel. tgadi@idc.ac.il.

1 Introduction

1.1 Motivation

We explore several classic problems in concurrent computing (election, mutual exclusion, consensus, and naming) when the number of processes which may participate is (denumerably) infinite. Partial information about the number of actually participating processes and the concurrency level is shown to affect the possibility and complexity of solving these problems. This paper surveys and generalizes work carried out in models with finite bounds on the number of processes, and proves several new results. These include improved bounds for election when participation is required (even for finitely many processes, as investigated by Styer and Peterson [SP89]) and a new adaptive starvation-free mutual exclusion algorithm for unbounded concurrency. We survey results in models with shared objects stronger than atomic registers, such as test&set bits, semaphores or read-modify-write registers, and update them for the infinite process case.

Processes: In most work on the design of shared memory algorithms, it is assumed that the number of processes is finite and *a priori* known. Here we investigate the the design of algorithms assuming no *a priori* bound on the number of processes. In particular, we assume that in an infinite run the number of active processes may be infinite. The primary motivation for such an investigation is to understand the limits of distributed computation—in particular, whether bounding the number of active processes is necessary in order to solve specific distributed problems. While in practice the number of processes will always be finite, algorithms designed for an infinite number of processes may scale well: their time complexity may depend on the actual contention and not on the total number of processes.

Concurrency: An important factor in designing algorithms where the number of processes is unknown, is the *concurrency level*, the maximum number of processes that may be active simultaneously, participating in the algorithm at the same instant of time. The actual concurrency in a given run is often called its *point contention*. We consider computability when the concurrency level imposes an *a priori* bound on the point contention. (A weaker notion of concurrency, *interval contention*, is not considered here.) We distinguish between the following concurrency levels:

- *finite*: There is a finite bound (denoted by c) on the maximum number of processes that are simultaneously active, over all runs. (The algorithms in this paper assume that c is known.)
- *bounded*: In each run, there is a finite bound on the maximum number of processes that are simultaneously active.
- *unbounded*: In each run, the number of processes that are simultaneously active in any state is finite but can grow without bound.

Participation: When assuming a fault-free model with *required participation* many problems are solvable using only constant space. In this model, every process must eventually execute its code. It had been considered, for example, for the symmetric, finite process model by Styer and Peterson [SP89]. However, a more interesting and practical situation is one in which participation is not required, as is more usually assumed when solving resource

allocation problems. For example, in the mutual exclusion problem a process can stay in the remainder region forever and is not required to try to enter its critical section.

We use the notation $[\ell, u]$ -*participation* to mean that at least ℓ and at most u processes participate. Thus, for a total of n processes (where n might be ∞) $[1, n]$ -participation is the same as saying that participation is not required, while $[n, n]$ -participation is the same as saying that participation is required. Requiring that all processes must participate does not mean that there must be a point at which they all participate at the same time. That is, the concurrency level might be smaller than the upper bound on the number participating processes. Notice also that if an algorithm is correct assuming $[\ell, u]$ -participation, then it is also correct assuming $[\ell', u']$ -participation, where $\ell \leq \ell' \leq u' \leq u$.

Thus, any solution assuming that participation is not required, is correct also for the case when participation is required, and hence it is expected that such solutions (for the case where participation is not required) may be less efficient and harder to construct.

1.2 Properties

We define two algorithm properties considered in the paper.

Adaptive algorithms: An algorithm is *adaptive* if the time complexity of processes' operations is bounded by a function of the actual concurrency. Time complexity is computed using a standard model, as was defined in [PF77], in which each primitive operation on an object is assumed to take no more than one time unit. In the case of mutual exclusion algorithms, we measure the maximum time between releasing the critical section, until the critical section is re-entered; this measure is called *system response time*. For models in which there is a minimum bound on the number of participating processes, we measure one time unit to the first point at which the minimum number of processes have begun executing the algorithm.

Symmetric algorithms: An algorithm is *symmetric* if the only way for distinguishing processes is by comparing (unique) process identifiers. In such algorithms process id's can only be written, read, and compared. In particular, identifiers cannot be used to index shared registers. Variants of symmetric algorithms can be defined depending on how much information can be derived from the comparison of two unequal identifiers. In this paper we assume that id's can only be compared for equality.

In particular, there is no total order on process id's in symmetric algorithms. Styer and Peterson discuss two types of symmetry [SP89]. Our notion corresponds to their stronger restriction, "symmetry with equality". Some of our asymmetric algorithms, presented later, satisfy their weaker symmetry restriction, "symmetry with arbitrary comparisons," in that they depend upon a total order of the process id's, but do not (for example) use id's to index arrays.

In Section 5.2 we study the *naming* problem in a still weaker model, in which processes are initially identical and utilize shared objects to acquire unique names.

1.3 Summary of results

We assume the reader is familiar with the definitions of the following problems:

1. The mutual exclusion problem, which is to design a protocol that guarantees mutually exclusive access to a critical section among a number of competing processes [Dij65];
2. The consensus problem, which is to design a protocol in which all correct processes reach a common decision based on their initial opinions [FLP85];
3. The (leader) election problem, which is to design a protocol where any fair run has a finite prefix in which all processes (correct or faulty) commit to some value in $\{0, 1\}$, and exactly one process (the leader) commits to 1. Although we require all processes to terminate in fair runs, we do not require the leader to be identified to the other processes. (In a model with infinitely many processes, identifying the leader obviously requires infinite space—the weaker $\{0, 1\}$ formulation makes the lower bounds more complex.) To prevent trivial solutions, we assume that the id’s of participating processes are not initially known. For asymmetric algorithms, we assume that process identities are natural numbers.
4. The *wait-free naming* problem which is to assign unique names to initially identical processes. Every participating process is able to get a unique name in a finite number of steps regardless of the behavior of other processes.

We show that even with a fixed bound on concurrency and required participation, election (using registers) requires infinite shared space. Among the novel algorithms presented are two demonstrating that either a single shared register of infinite size, or infinitely many shared bits, suffice for both election and consensus. (In addition, the first algorithm is adaptive.) If in addition test&set bits are used, then solving the above problem requires only finite space; however, a result of Peterson ([Pet94]) implies that the more complex problem of starvation-free mutual exclusion (bounded concurrency, participation not required) still requires infinite space. In fact, even using read-modify-write registers to solve this problem, a result of Fischer *et al* ([F⁺89]) implies that infinite space is required. However, Friedberg and Peterson ([FP87]) have shown that using objects such as semaphores that allow waiting enables a solution with only constant space.

When there is no upper bound on the concurrency level, we show that using an infinite number of registers is a necessary and sufficient condition for solving the mutual exclusion problem. (The algorithm presented is adaptive, symmetric, and satisfies starvation-freedom.) We then show that even infinitely many test&set bits do not suffice to solve (wait-free) naming assuming unbounded concurrency. However, naming can be solved assuming bounded concurrency using test&set bits only, hence it separates bounded from unbounded concurrency.

The tables below summarize the results discussed in this paper. As indicated, many are derived from existing results, generally proven for models where the number of processes is finite and known. We use the following abbreviations: DF for deadlock-free; SF for starvation-free; mutex for mutual exclusion; U for Upper bound; L for Lower bound; RW for atomic read/write registers; T&S for test&set bits; wPV for weak semaphores; and RMW for read-modify-write registers. (The default is “No” for the adaptive and symmetric columns. All lower bounds hold for the most general, non-adaptive and asymmetric case.)

Problem		Model			Result					
Name	Bound	Concurrence $c > 1$	Participation	Properties		Space		Thm #	Using results from	
				adaptive?	sym-metric?	#	size			
1	Election (Upper bounds hold also consensus.)	L ₁	c	$[c, \infty]$			∞ state space		2.1	
		U	c	$[c, \infty]$			∞	2	2.3	
		U	2	$[2, \infty]$	Y	Y	1	∞	2.4	
		U	c	$[c, \infty]$	Y	Y	2	∞	2.5	
		U	c	$[c, \infty]$	Y		1	∞	2.6	
		U	c	$[1, \infty]$		Y	$\log c + 1$	∞	3.1	[SP89]
		L ₂	c	$[1, \infty]$			$\log c + 1$	2	3.1	[SP89]
2	DF mutex	U	c	$[1, \infty]$		Y	c	∞	3.1	[SP89]
		L	c	$[1, \infty]$			c	2	3.1	[BL93]
3	Election SF mutex	L	bounded	$[1, \infty]$			∞	2	3.1	1L ₂ ,2L
		U ₁	unbounded	$[1, \infty]$	Y	Y	∞	∞	3.2	
		U ₂	unbounded	$[1, \infty]$			∞	2	3.7	
4	Implement Test & Set	L	bounded	$[1, \infty]$			∞	2	4.1	3L,5U
		U	unbounded	$[1, \infty]$	Y	Y	∞	∞	4.1	3U ₁
		U	unbounded	$[1, \infty]$			∞	2	4.1	3U ₂

Results using atomic registers in a fault-free model

Problem		Model					Result					
Name	Bound	Concurrence $c > 1$	Participation	Fault tolerance	Shared objects	Properties		Space		Thm #	Using results from	
						adaptive?	sym-metric?	#	size			
5	Election DF mutex	L	c	$[\infty, \infty]$	0	T&S			1	2		trivial
		U	unbounded	$[1, \infty]$	∞	T&S	Y	Y	1	2		trivial
6	SF mutex	L	bounded	$[1, \infty]$	0	RW,T&S			∞	2	4.2	[Pet94]
		U	unbounded	$[1, \infty]$	0	RW			∞	2	3.7	3U ₂
7	Naming	U	bounded	$[1, \infty]$	∞	T&S	Y	Y	∞	2	4.3	[AT96] [MA95]
		L	bounded	$[1, \infty]$	∞	T&S	Y	Y	∞	2	4.3	
		L	unbounded	$[1, \infty]$	∞	T&S			no alg		4.3	
8	Consensus	L	2	$[\infty, \infty]$	1	RMW			2	2		trivial
		U	unbounded	$[1, \infty]$	∞	RMW	Y	Y	1	3		trivial
9	SF mutex	L	bounded	$[1, \infty]$	0	RMW			1	∞	5.1	[F ⁺ 89]
		U	unbounded	$[1, \infty]$	0	RMW	Y	Y	1	∞	5.1	[B ⁺ 82]
		U	unbounded	$[1, \infty]$	0	RW			∞	2	3.7	3U ₂
10	SF mutex	U	unbounded	$[1, \infty]$	0	RW,wPV	Y	Y	2,2	2	5.2	[FP87]

Results using stronger shared objects

1.4 Related Work

In this paper we study the problem of computing with infinitely many processes in models with communication primitives stronger than atomic read/write registers and studying problems such as mutual exclusion that admit fault-free solutions. However, we do not consider wait-free algorithms using atomic registers only. An earlier version of this paper

appeared in [MT2000].

In [GMT2001], wait-free computation using only atomic registers is considered. It is shown that bounding concurrency reveals a strict hierarchy of computational models, of which unbounded concurrency is the weakest model. Nevertheless, it is demonstrated that adaptive versions of many interesting problems (collect, snapshot, renaming) are solvable even in the unbounded concurrency model. Wait-free randomized consensus algorithms for infinitely many processes using only atomic registers has been explored in [ASS2002], where it is also observed that standard universal constructions based on consensus continue to work with infinitely many processes with only slight modifications. In [MT2003], a general construction is presented, which implements wait-free consensus for infinitely many processes from any known solution for wait-free consensus (deterministic or randomized) for only finitely many processes. In [CM2002], the active disk paxos protocol is implemented for infinitely many processes, which facilitates a solution to the consensus problem with an unbounded number of processes.

A starvation-free mutual exclusion algorithm has long been known for unbounded concurrency, using two registers and two weak semaphores (see Theorem 5.2) [FP87]. In [BMT95], various types of shared counters are implemented in a model where nothing can be assumed in advance about the number or the identities of the processes that might access the counters. Wait-free solvability of tasks when there is no upper bound on the number of participating processes has also been investigated [GK98], but in this earlier work no run has an infinite number of participating processes (see also [Gafni2002]). An underlying, automata-theoretic formulation of unbounded concurrency may be found in e.g. [L⁺93, Lyn96].

All previously published mutual exclusion algorithms using atomic registers assume that the number of processes are finite and known. The question whether there exists an adaptive mutual exclusion algorithm using atomic registers was first raised in [MT93], where an adaptive algorithm is presented for a given working system, which is useful provided process creation and deletions are rare. In [MT93], the term *contention sensitive* was used but later the term *adaptive* became commonly used. In [AT92], it is proven that no adaptive algorithm exists when time is measured by counting all accesses to shared registers. Hence we use another time complexity measure (defined previously), called system response time [PF77], for evaluating the complexity of adaptive mutual exclusion algorithms. Other complexity measures were also considered in the literature, for example counting only the number of remote memory references [AK2000], or assuming that busy-waiting is counted as a single step [AB2002, AST99].

In [CS94], the only previously known adaptive mutual exclusion algorithm was presented, in a model where it is assumed that the number of processes (and hence concurrency) is finite. The algorithm exploits this assumption to work in bounded space, and does not work assuming unbounded concurrency. In [Lam87], a fast algorithm is presented which provides fast access in the absence of contention. However, in the presence of any contention, the winning process may have to check the status of all other n processes (i.e. access n different shared registers) before it is allowed to enter its critical section. Recent results about adaptive mutual exclusion algorithms for finite and a priori known number of processes include [AB2002, AK2000, AST99, AST2002, CS94, Tau2004].

A symmetric election algorithm is presented in [SP89], for n processes (with concurrency level n), which uses only three atomic registers. Several other related references are mentioned in context later in this paper.

2 Atomic registers: participation is required

This section demonstrates that requiring a minimum number of processes to participate (the required participation model) is a powerful enabling assumption: the problems we study are solvable with a small number of registers. However, reducing the number of registers does not mean that it is also possible to achieve a finite shared state space. We begin by showing that any solution to election in this model requires infinite shared space. Then we show that election can be solved either by using an infinite number of binary registers, or a single register of infinite size.

2.1 A lower bound for election when the concurrency is $c > 1$

Even when participation is required, (and equal to the concurrency) any solution to election for infinitely many processes must use infinite shared space. This holds even if processes do not need to learn the identity of the leader as part of the election algorithm.

Theorem 2.1 *There is no solution to election with finite concurrency $c > 1$ and with $[c, \infty]$ -participation using finite shared memory (finitely many registers of finite size).*

Proof: Assume to the contrary that there is such an algorithm using finite space. Consider the solo run $solo(p)$ of each process p from the initial state, and define $write(p)$ to be the sequence of states of the shared memory in $solo(p)$. (Because $c > 1$, $solo(p)$ may be infinite, as process p waits for the arrival of additional processes.) For each process p , let $repeat(p)$ be the first state of the shared memory that repeats infinitely often in $solo(p)$. (Recall that the algorithm uses finite space by assumption. Also, if $write(p)$ is finite, define $repeat(p)$ to be the state after the last write in $write(p)$.)

Let $\beta(p)$ be the finite prefix of $write(p)$ that precedes the first instance of $repeat(p)$. For each such $\beta(p)$, we construct a *signature* subsequence, $sign(p)$, by removing repetitive states and the intervening steps, as follows: Let $\beta(p) = x_0, x_1, x_2, \dots$, where the x_i are the successive states of the shared memory in $\beta(p)$. Suppose x_j is the first state that repeats in $\beta(p)$, and that x_k is the last state in $\beta(p)$ such that $x_j = x_k$. Remove the subsequence x_{j+1}, \dots, x_k from $\beta(p)$. The resulting sequence $x_0, \dots, x_j, x_{k+1}, \dots$ is a subsequence of $\beta(p)$ with strictly fewer repeating states. Repeat this step (finitely many times) until no state repeats—the resulting sequence is (the signature) $sign(p)$.

Since there are only finitely many states of the shared memory, there exists a single state s and a single sequence γ such that $s = repeat(p)$ and $\gamma = sign(p)$ for infinitely many processes, p . Let $C = \{p_1, \dots, p_c\}$ be any subset of c of these processes.

Lemma 2.2 *The solo runs of processes in C are compatible: There exists a single infinite run $elect(C)$ of the processes in C in which each process p_i takes the same steps as in $solo(p)$.*

Proof: The run $elect(C)$ is constructed as follows: Let $\gamma = y_0 y_1, \dots, y_k$. Run p_1 until its first write to the shared memory that changes the shared memory state. Do the same for p_2 through p_c . Now each process p_i is about to execute its first write that actually changes the value of shared memory, from y_0 to y_1 . Run process p_1 until it takes the last step in $\beta(p_1)$ that leaves the memory in state y_1 . (Note again that y_0 and y_1 differ by the value of exactly one word of the shared memory, which is exactly the word that each of p_2 through p_c are

preparing to change, by writing the same value, y_1 .) Repeat for p_2 through p_c , so that each process p_i in $\{p_1, \dots, p_c\}$ has executed the steps of $solo(p_i)$ up to the last occurrence of y_1 in $\beta(p)$. Now repeat this loop $k - 1$ more times, for each remaining state in γ , so that in the j 'th iteration, each process p_i in $\{p_1, \dots, p_c\}$ has executed the steps of $solo(p_i)$ up to the last occurrence of y_j in $\beta(p)$.

The resulting finite run leaves the shared memory in state s , and each of the c processes p_i has taken all the steps in $\beta(p_i)$. Since the state s repeats infinitely often in the remainder of each of the solo runs $solo(p_i)$, the run can be extended, using a round-robin schedule, appending a finite sequence of steps of each process in turn, always returning the memory to state s . The resulting run is $elect(C)$. ■

To conclude the proof of Theorem 2.1, note that since $elect(C)$ has concurrency at most c and participation c , one of the participants, p_i , must be elected. Similarly, for a set of processes $C' = p'_1, \dots, p'_c$, C' disjoint from C , in the run $elect(C')$ some process $p_{i'}$ is elected. Finally, including both p_i and $p_{i'}$ in a subset C'' of $C \cup C'$ of size c results in a run $elect(C'')$ in which both p_i and $p_{i'}$ are elected, contradiction. ■

2.2 Algorithms for election and consensus for concurrency c

The previous theorem shows that election for infinitely many processes requires unbounded shared memory. In this section, we present algorithms for the required participation case that bound either the number of shared words or their size. We first study the scenario in which the concurrency is equal to participation (concurrency c and participation $[c, \infty)$). We show that election can be solved either by (1) using an infinite number of atomic bits, or (2) using a single register of infinite size. We also present two simple symmetric algorithms.

Theorem 2.3 *For any finite concurrency c , with $[c, \infty)$ -participation, there are non-adaptive asymmetric solutions to election (and consensus) using an infinite number of atomic bits.*

Proof: We identify each process with a unique natural number, and assume 0 is not assigned to any process. The algorithm uses an infinite array $b[0], b[1], \dots$ of bits, which are initially 0.

The first step of process i is to read $b[0]$. If the value seen is 1, process i knows that a leader has already been elected and terminates. Otherwise, process i sets $b[i]$ to 1. Then, process i scans the array until it notices that c other bits (other than $b[0]$) are set. In order not to miss any bit that is set, it scans the array according to an infinitely repeating enumeration—a schedule that visits each bit an infinite number of times. (A canonical example is to first read $b[1]$; then $b[1], b[2]$; then $b[1], b[2], b[3]$, etc.) Once a process notices that c bits are set, it sets $b[0]$ to 1, chooses the process with the smallest id among the c that are set to 1 as the leader, and terminates. (By scanning the bits after reading $b[0]$ as 1, the other processes can also learn the id of the leader.)

Only processes which read $b[0] = 0$ will write any bits, and no such process terminates until it sees c other bits set. Hence, the first c processes do not terminate until all of them read $b[0] = 0$ and set their bits. The first one to terminate sets $b[0]$ to its final value of 1. Only then can the $c + 1$ 'st process become active. Hence, exactly $c + 1$ bits will eventually be set.

Process i 's program

Shared: $(Leader, Marked)$: (Process id, boolean) initially $(0, 0)$

Local: $local_leader$: Process id

```

1  if  $Marked = 0$  then
2       $(Leader, Marked) := (i, 0)$ 
3      await  $(Leader \neq i)$  or  $(Marked = 1)$ 
4       $local\_leader := Leader$ 
5       $(Leader, Marked) := (local\_leader, 1)$ 
6  fi
7  return $(Leader)$ 

```

Figure 1: Symmetric election for concurrency 2

This solution for election can be trivially modified, using one other bit, $b[-1]$, to solve consensus. Each process i (of the first c) which reads $b[0] = 0$ sets $b[-1]$ to its input value before setting $b[i] = 1$ and continuing with the rest of the algorithm. Before terminating, each process reads and returns the value of $b[-1]$. The consensus value will be the input of the last, of the first c processes participating, to write to $b[-1]$. ■

A symmetric election algorithm is presented in [SP89], for n processes (with concurrency level n), which use only three atomic registers. Below we present three election algorithms: (1) a symmetric algorithm for concurrency 2 using one register, (2) a symmetric algorithm, using an idea from [AGM98], for concurrency c using two registers, and finally, (3) an asymmetric algorithm for concurrency c using only a single register.

Theorem 2.4 *For finite concurrency $c = 2$ with $[2, \infty]$ -participation, there are adaptive symmetric solutions to election and consensus using one atomic register.*

Proof: The algorithm in Figure 1 is very simple. The register value contains two fields, $Leader$ and $Marked$, respectively, where $Leader$ is a process id (or 0) and $Marked$ is boolean. The statement **await condition** is used as an abbreviation for **while** \neg *condition* **do skip**, and hence, may involve many accesses to the shared memory.

In this very simple algorithm the second process to write to the shared register is elected. The first process to write spins until the second process writes, then marks the second process as the winner. The concurrency bound guarantees no other process will interfere.

This election algorithm can be easily converted into a consensus algorithm by appending the input value (0 or 1) to each process id. The input value of the leader is the consensus value. ■

The next algorithm, solving asymmetric election for concurrency $c \geq 2$, is similar to the previous, in that the last participant to write the $Leader$ field is then $Marked$ as the leader. A second register, $Union$, is used by the first c participants to ensure that they have all finished writing their id into $Leader$. (Note, that in doing so, they learn each other's id's.) Because c is also the number of required participants, each of the first c participants can spin on the $Union$ register until the number of id's in it is c .

Process i 's program.

Shared:

$(Leader, Marked)$: (Process id, boolean), initially $(0, 0)$

$Union$: set of at most c process id's, initially \emptyset

Local:

$local_leader$: Process id

$local_union1$: set of at most c process id's

$local_union2$: set of at most c process id's, initially $\{i\}$

```
1  if  $Marked = 0$  then  $(Leader, Marked) := (i, 0)$  fi
2   $local\_union1 := Union$ 
3  while  $(|local\_union1| < c) \wedge (Marked = 0)$  do
4      if  $\neg(local\_union2 \subseteq local\_union1)$  then
5           $Union := local\_union1 \cup local\_union2$ 
6      fi
7       $local\_union2 := local\_union1 \cup local\_union2$ 
8       $local\_union1 := Union$ 
9  od
10  $local\_leader := Leader$ 
11  $(Leader, Marked) := (local\_leader, 1)$ 
12 return $(Leader)$ 
```

Figure 2: Symmetric election for concurrency c

Theorem 2.5 *For any finite concurrency c with $[c, \infty]$ -participation, there are adaptive symmetric solutions to election and consensus using two atomic registers.*

Proof: The algorithm in Figure 2 uses a key idea used previously in a symmetric algorithm for a general function evaluation problem, due to Attiya, *et al* [AGM98]. As in the Attiya, *et al* algorithm, the first register, $Union$, holds a set of at most c processes id's. Processes engage in a phase of reading the current set of values, comparing it to the set they have already seen, and writing when they know a value not reflected in the set read. This phase ends when a process reads c different values. (In the Attiya, *et al* algorithm, this register holds a set of inputs to a collectively-computed function. This phase ends when a process observes enough inputs to uniquely determine the function value, which is written in the second register, providing the termination condition and output value for all participating processes.)

In the algorithm in Figure 2, the second register is used both to “gate” access to the $Union$ register, and to break symmetry on the set of process id's. It contains two fields: a process id (in the $Leader$ field) and a bit (the $Marked$ field). Processes access the $Union$ field only if they read the $Marked$ field as 0. In addition, no process terminates until it sets the $Marked$ field to 1. The concurrency bound guarantees that the $c + 1$ 'st process cannot take a step until one of the first c processes terminate, by which time the $Marked$ field is set. Hence, only the first c processes access the $Union$ register. We prove below that one of these processes must eventually see c different process id's (and go on to set $Mark$). Hence, the leader will be the c 'th process to write to $Leader$.

The proof begins by arguing that all the processes that terminate elect exactly the same participant as a leader. Note first that only participating process id's can appear

in the shared *Union* register. A process sets *Marked* from 0 to 1 and terminates only after it determines that c processes are already participating. Thus, the first c processes to participate will first see that *Marked* is not set, and no later process can participate until one of these terminates, after setting *Marked*. Hence, only the first c processes will write their id's into *Leader*. Moreover, no process sets *Marked* until the last of these c processes have written to *Leader*. Hence, any processes that terminate will elect the last process to write its id to *Leader*.

To prove termination, consider the finite partial order defined by the subset relation on the subsets of the set of names C of the first c processes. We can consider each of the first c processes as moving from one point to another in this partial order as their *local_union2* variables are updated. (Several processes may move to the same point.) Because these variables are only updated by union with their previous value, processes can only move upward in this order. (Call the set of id's in *local_union2* the set that the associated process currently *knows about*.) As long as *Marked* is not set, we argue that the only stable state of this marking of the partial order is one in which every process knows about every other. (That is, $|local_union2| = c$ for each of the first c processes.)

Consider any other reachable state, and assume that the marking of the partial order is stable in a run extending it, in which each of the first c processes either takes an infinite number of steps or terminates. This state contains minimal elements: processes with $|local_union2| < c$ and such that no other process knows only about a strict subset of *local_union2*. Once a minimal process writes its value of *local_union2* to *Union*, no strict subset of that set can be written. (By monotonicity of the *local_union2* variables.) Each such process can write at most once in the run extending this state, as it writes only if it learns of a new process, which by assumption it cannot. So in its (infinitely many) subsequent reads, it must see exactly the set it currently knows. But eventually a process that knows a different id will write to the *Union* register and then be read by a minimal process, contradiction.

It follows that at some point a process will learn that c processes are participating (i.e., $Union = C$ and $|Union| = c$) and set *Marked* to 1, allowing the rest of the processing to terminate.

Establishing adaptivity means bounding the time complexity of each participating process by a function of c . For a subset K of C , define $W(K)$ as the maximum number of times any single process can write K to *Union*. (Note that only processes in K can write K to *Union*.) Each write must be preceded by the read of a strictly smaller set. Hence, for any $p \in C$, we have $W(\{p\}) = 1$, and $W(K) \leq \sum_{K' \subset K} |K'| \cdot W(K')$. Hence the number of distinct writes to *Union* is bounded by $\sum_{K \subset C} |K| \cdot W(K)$, a function of c .

This completes the correctness proof of the election algorithm.

As above, this election algorithm can be easily converted into a consensus algorithm, by adding the input value (0 or 1) to each process id. ■

Note that the termination argument above holds even if some processes stop updating the *Union* register, as long as the set of active processes knows about everyone. This observation is important to the next algorithm, an asymmetric solution to election (and consensus) using only a single atomic register. Asymmetric algorithms allow comparisons between process id's (we assume a total order). However, since the identities of the participating process are not initially known to all processes, the trivial solution where all processes choose process 1 as a leader is excluded.

Theorem 2.6 *For any finite concurrency c and $[c, \infty]$ -participation there is an adaptive asymmetric solution to election and consensus using one atomic register.*

Proof: The one-register asymmetric algorithm in Figure 3 can be understood as an adaptation of the previous, two-register algorithm for the symmetric case (Theorem 2.5). As in an asymmetric one-register algorithm due to Attiya, *et al* [AGM98], it modifies the symmetric algorithm by adding a leader-lead handshake phase. (The algorithm here has slightly different termination conditions than the function evaluation algorithm, to prevent late processes from running concurrently with the first c active processes. We also prove adaptivity of the algorithm.)

As before, each of the first c processes (denoted by the set C) repeatedly updates the shared *Union* register with any new information it obtains. The process elected as leader is the process with the minimum id among the first c processes to participate in the algorithm. Termination detection among the first c processes, is accomplished by establishing an explicit handshake between the leader and each of the other $c - 1$ initial processes.

In outline, once a process p observes the size of the set in *Union* to be at least c (we argue below exactly c), it knows the identity of the leader (the minimum element in this set). If p is not in the set (is not one of the first c participants) it simply terminates. Otherwise, p performs a handshake with the leader, and terminates once it observes that all the $c - 1$ handshakes have completed. The handshake is carried out using additional fields in the register: a counter, *Sig*, by which the leader signals to the initiating process with rank *Sig* in *Union*, and a bit, *Ack*, by which this process acknowledges. That is, process p waits until the leader sets *Sig* to p 's rank among the c initial processes, and then sets *Ack* to 0 to acknowledge the signal. Except for acknowledging the signal, p stops participating in the updates to *Union*, but waits to terminate until the final signal value of c is posted. (This precludes any but the first c active processes from participating until the register value has the stable value of $(C, c, 1)$. Moreover, at the point that process p acknowledges the signal, it knows that the leader knows all c initial id's, and will eventually refresh *Union* if necessary.) The leader itself initiates each handshake in order once it knows the id's of all c initial processes. Slow processes may overwrite handshakes in progress, which the leader (or higher-ranking processes) will detect and refresh.

As in the previous algorithm, no process exits the while loop in line 2 until c processes have entered the loop. We argue below that none of the initial processes terminates until the shared register contains the final and stable value $(C, c, 1)$. It follows that all processes that terminate elect the same leader.

To argue that every process terminates, the main difficulty is arguing that the initial c processes terminate. Note that no process stops updating *Union* until it knows the id's of c participating processes, (the initial while loop in line 2) and indeed, until it also knows that the leader knows the id's of c participating processes (the update in line 16 that is skipped once the leader's signal is acknowledged). As in the argument in the preceding algorithm, this implies that all processes in C will eventually learn the id's in C .

By induction on the rank $r \geq 1$ of process p in C , process p eventually reads its signal from the leader, the leader eventually reads its acknowledgment, and (using the explicit *acked* flag to avoid unnecessary refreshes of the register) process p never executes another write operation. Hence, the leader will eventually set the register to $(C, c, 1)$ and terminate. The other of the initial participants (who have stopped writing but are still reading) will eventually read *Sig* = c and terminate, and as we argued above, the later processes will all see c different values in *Union* and terminate immediately.

Process i 's program

Shared: $(Union, Sig, Ack)$: $Union$ a set of at most c process id's,
 $0 \leq Sig \leq c, Ack \in \{0, 1\}$, initially $(\emptyset, 0, 0)$

Local: $local_union1$: set of at most c process id's
 $local_union2$: set of at most c process id's, initially $\{i\}$
 $myrank, local_sig1, local_sig2$: integers between 0 and c , initially 0
 $local_ack1, acked$: Boolean, initially 0
 $leader$: process id

```

1  ( $local\_union1, local\_sig1, local\_ack1$ ) := ( $Union, Sig, Ack$ )
2  while  $|local\_union1| < c$  do
3      if  $\neg(local\_union2 \subseteq local\_union1)$  then
4          ( $Union, Sig, Ack$ ) := ( $local\_union1 \cup local\_union2, 0, 0$ )
5      fi
6       $local\_union2 := local\_union1 \cup local\_union2$ 
7      ( $local\_union1, local\_sig1, local\_ack1$ ) := ( $Union, Sig, Ack$ )
8  od
9   $local\_union2 := local\_union1$ 
10  $leader := \min\{q : q \in local\_union1\}$ 
11 if  $i \notin local\_union1$  then
12     return( $leader$ )
13 elseif  $i \neq leader$  then
14      $myrank := |\{h \in local\_union1 : h < i\}|$ 
15     while  $local\_sig1 < c$  do
16         if  $(|local\_union1| < c) \wedge (acked = 0)$  then /* register is stale, haven't acked */
17             ( $Union, Sig, Ack$ ) := ( $local\_union2, 0, 0$ )
18         elseif  $(local\_sig1 = myrank) \wedge (local\_ack1 = 1)$  then
19             ( $Union, Sig, Ack$ ) := ( $local\_union2, myrank, 0$ ) /* acknowledge signal */
20              $acked := 1$ 
21         fi
22         ( $local\_union1, local\_sig1, local\_ack1$ ) := ( $Union, Sig, Ack$ )
23     od
24     return( $leader$ )
25 else /*  $|local\_union2| = c, i = leader$  */
26     while  $local\_sig1 < c$  do
27         if  $(|local\_union1| < c) \vee (local\_sig1 < local\_sig2)$  then
28             ( $Union, Sig, Ack$ ) := ( $local\_union2, local\_sig2, local\_ack2$ ) /* register is stale */
29         elseif  $local\_ack1 = 0$  then /* signal acknowledged */
30              $local\_sig2 := local\_sig2 + 1$ 
31             ( $Union, Sig, Ack$ ) := ( $local\_union2, local\_sig2, 1$ ) /* send new signal */
32         fi
33         ( $local\_union1, local\_sig1, local\_ack1$ ) := ( $Union, Sig, Ack$ )
34     od
35     return( $i$ )

```

Figure 3: (Asymmetric) election for concurrency c

As in the previous algorithm, establishing adaptivity means bounding the time complexity of each participating process by a function of c . The proof of Theorem 2.5 defined a recurrence based on the finite partial order of the subsets of the set C of the first c active processes. To prove the algorithm in Figure 3 adaptive, extend the subset partial order with the $2c - 1$ handshake values in the linear order of the handshake:

$$C < (C, 1, 1) < (C, 1, 0) < (C, 2, 1) < (C, 2, 0) < \dots < (C, c, 1)$$

The same recurrence argument as in the proof of Theorem 2.5 establishes adaptivity by bounding the number of times any value can be written. ■

3 Atomic registers: participation is not required

As mentioned above, any problem solution in a model in which participation is not required, is also correct when participation is required, and hence such solutions may be more difficult to construct. The main new result of this section is an adaptive starvation-free mutual exclusion algorithm for unbounded concurrency.

3.1 Consensus, election and mutual exclusion for concurrency c

We next consider the number and size of registers required to solve consensus, election and mutual exclusion in a model in which participation is not required (i.e., $[1, \infty]$ -participation). Earlier work on these problems for a finite number of processes establishes some bounds:

Theorem 3.1 *For concurrency level c , the number of atomic registers that are:*

- (1) *necessary [BL93] and sufficient [SP89] for solving deadlock-free mutual exclusion, is c ;*
- (2) *necessary and sufficient for solving election, is $\log c + 1$ [SP89];*
- (3) *sufficient for solving consensus, is $\log c + 1$ [SP89].*

Proof:

(1) Any deadlock-free mutual exclusion algorithm for n processes must use at least n shared registers [BL93]: this bound applies immediately to concurrency.

There is a symmetric deadlock-free mutual exclusion algorithm for n processes which uses only n shared registers [SP89]. This algorithm works without modification for an infinite number of processes, provided that the concurrency level is no more than n .

(2) Any election algorithm for n processes must use at least $\log n + 1$ shared registers [SP89]: as above, this bound applies immediately to concurrency.

There is a symmetric election algorithm for n processes which uses only $\log n + 1$ shared registers [SP89]: it also works without modification for an infinite number of processes, provided that the concurrency level is no more than n .

(3) The symmetric election algorithm in [SP89] can be easily modified to solve consensus with the same number of registers: Each process appends its vote (0 or 1) to its id. Processes choose the leader's vote as the consensus value. ■

3.2 Mutual exclusion algorithms for unbounded concurrency

Theorem 3.1 implies that when the concurrency level is not finite, an infinite number of registers are necessary for solving the election and mutual exclusion problems. But do an infinite number of registers suffice for solving mutual exclusion for unbounded concurrency, when participation is not required? (We observe, based on a result of Yang and Anderson [AY96], that if we restrict the number of processes that can try to write the same register at the same time, the answer is no.) We present two algorithms that answer this question affirmatively. The first is an adaptive and symmetric mutual exclusion algorithm using infinitely many infinite-sized registers. The second is neither adaptive nor symmetric, but uses only (infinitely many) bits.

Theorem 3.2 *There is an adaptive symmetric solution to starvation-free mutual exclusion for unbounded concurrency (and hence also to election and consensus) using an infinite number of registers.*

These problems can all be solved by simple adaptations to the deadlock-free mutual exclusion algorithm presented in Figure 6 below. This algorithm has the following interesting properties:

1. It works for unbounded concurrency.
2. In the absence of contention, only *eight* accesses to the shared memory are needed (*seven* in the entry code and *one* in the exit code).
3. It is adaptive – its time complexity (even in the worst case, measured from exit of the critical section by a process to the next entrance by any process) is a function of the actual number of contending processes.
4. It is symmetric – identifiers are only written, read and compared for equality, but are not ordered or used to index shared registers.

Except for the non-adaptive algorithm in the next subsection, we know of no mutual exclusion algorithm (using atomic registers) satisfying the first property. The algorithm is defined formally in Figure 6, and is built out of simple building blocks called *splitters* which are introduced next.

3.2.1 Splitters

In [Lam87], Lamport presented a fast mutual exclusion algorithm which provides fast access in the absence of contention. As emphasized by Moir and Anderson [MA95], his algorithm makes use of a shared object called a *splitter*. This object is shown in Figure 4. Each process that invokes the splitter moves either *down*, *right* or *wins*. In any execution, define the *latecomers* to be the processes that invoke the splitter after the first process exits it. Let n be the number of *early* processes that invoke the splitter before this first process exits it. Properties of Lamport’s splitter are listed below.

Proposition 3.3 *In any execution of Lamport’s splitter, the following properties hold:*

1. *At most $n - 1$ early processes move right,*

2. at most $n - 1$ processes move down, or at least one latecomer moves right.
3. at most one process wins, and if $n = 1$ then exactly one process wins,
4. the latecomers all move right, and
5. the splitter is wait-free.

Process i 's program

Shared:

x : integer (the initial value is immaterial)

y : boolean, initially 0

```

1  $x := i$ 
2 if  $y = 1$  then goto right fi
3  $y := 1$ 
4 if  $x \neq i$  then goto down
5 else goto win fi

```

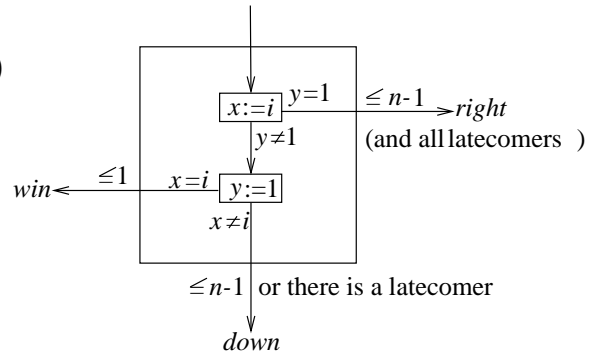


Figure 4: Lamport's splitter.

Proof: To see, for example, why the third property holds, assume to the contrary that two processes i and j both *win*. Assume without loss of generality that process i tests the value of x at statement 4 after j does so. This implies that x is not written by any process between i 's assignment $x := i$ in statement 1 and i 's read of x in statement 4. Thus, j read of x in statement 4 preceded i 's assignment in statement 1, which in turn implies that j assigned 1 to y in statement 3 before i 's read of y in statement 2. Thus, i must have read $y = 1$ at statement 2 and then “moved right”, a contradiction. Similar arguments establish the other properties. ■

Next, we modify Lamport's splitter, exchanging wait-freedom for starvation-freedom, but obtaining a new safety property:

Lemma 3.4 *The new splitter in Figure 5 satisfies the first four properties of Lamport's splitter (as listed in Proposition 3.3), is starvation-free, and satisfies the following additional property: If a process wins then nobody moves down.*

Proof: Starvation-freedom follows because Lamport's splitter is wait-free and the conditions satisfying the *await* statement are stable. The second property of Lemma 3.3 assures that some process must exit Lamport's splitter to the *right* or by *winning*, and perform the assignment enabling all *await* statements to terminate. To see why the new safety property holds, observe that:

1. a process that reaches the *await* statement in line 4, (exiting Lamport's splitter going down) can move *down* in the new splitter only if b is set to 1 *before* z is set to 1, and
2. a process can *win* in the new splitter only if b is set to 1 *after* z is set to 1.

Thus, if a process *wins* in the new splitter, every process that reaches the *await* statement in line 4 moves *right* in the new splitter. The winning process is also the (unique) process to *win* in Lamport's splitter, and so nobody moves *down*.

The other four safety properties of Lamport’s splitter also hold for the new splitter. Notice first that any process that is an *early* process in the embedded Lamport’s splitter is also an *early process* in the new splitter, and any process that is a *latecomer* in new splitter is a *latecomer* in the embedded Lamport’s splitter. The third and fourth safety properties follow immediately from the corresponding conditions in Lemma 3.3. To see the first safety property holds, observe that one of the processes that are *early* in Lamport’s splitter exits it as a *winner* or *down*. In the first case, we are done, so assume no process *wins* in Lamport’s splitter. Then the processes exiting Lamport’s splitter *down* will continue and exit the new splitter going *down*.

Finally, we prove that the second safety property holds. If some process moves *right* in the new splitter (either an *early* process or a *latecomer*) then we are done. If no process moves *right* in the new splitter, then also no process moves *right* in Lamport’s splitter, and by the second safety property in Lemma 3.3, exactly one process (*early* in both splitters) *wins* in Lamport’s splitter, and then in the new splitter. (Since b is never set to 1.) ■

Process i ’s program

Shared:

x : integer (the initial value is immaterial)

b, y, z : boolean, initially 0

```

1   $x := i$ 
2  if  $y = 1$  then  $b := 1$  goto right fi
3   $y := 1$ 
4  if  $x \neq i$  then await  $((b = 1) \text{ or } (z = 1))$ 
5      if  $z = 1$  then goto right
6      else goto down fi
7  else  $z := 1$ 
8      if  $b = 0$  then goto win
9      else goto down fi fi

```

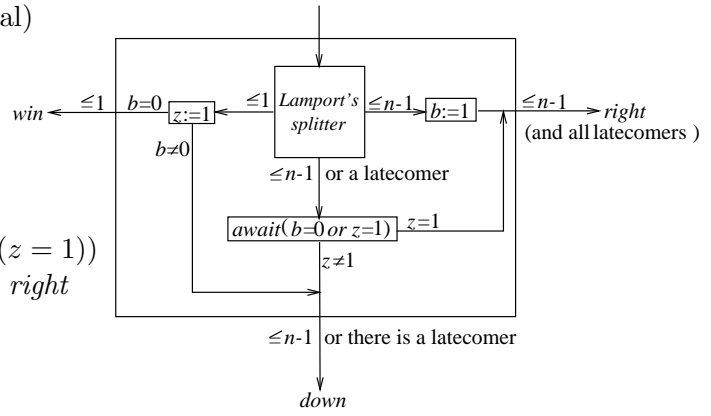


Figure 5: The code of the *new splitter*.

3.2.2 An adaptive mutual exclusion algorithm

Using the properties of new splitter, it is possible to solve mutual exclusion by interconnecting a collection of splitters in an (infinite) chain, so that processes that move *down* enter the next splitter in the chain, as is illustrated in Figure 6:

Lemma 3.5 *The algorithm in Figure 6 is an adaptive, deadlock-free mutual exclusion algorithm for unbounded concurrency.*

Proof: In this solution, the processes compete in levels, where each level is implemented as a separate new splitter. Each level either has a *winner*, or eliminates at least one competing process (those that move *right*). Since all the *latecomers* move *right* at the first level, only finitely many processes can move *down* to subsequent levels. Hence, within finitely many levels, there is either a *winner* or a single process moves down—but this process will *win* this final level. The *winner* enters its critical section, and in its exit code it publishes the index

Process i 's program

Shared:

$global_level$: integer, initially 0

$x[0..\infty]$: array of integers (the initial values are immaterial)

$b[0..\infty], y[0..\infty], z[0..\infty]$: array of boolean, initially all 0

Local:

my_level : integer (the initial value are immaterial)

```

0    $my\_level := global\_level$ 
1    $start: x[my\_level] := i$ 
2   if  $y[my\_level]$  then  $b[my\_level] := 1$  goto right fi
3    $y[my\_level] := 1$ 
4   if  $x[my\_level] \neq i$  then await  $((b[my\_level] = 1) \text{ or } (z[my\_level] = 1))$ 
5       if  $z[my\_level] = 1$  then  $goto\ right$ 
6       else  $goto\ down$  fi
7   else  $z[my\_level] := 1$ 
8       if  $b[my\_level] = 0$  then  $goto\ win$ 
9       else  $goto\ down$  fi fi

10   $right: await\ my\_level < global\_level$ 
11   $my\_level := global\_level$ 
12  goto  $start$ 

13   $down: my\_level := my\_level + 1$ 
14  goto  $start$ 

15   $win: critical\ section$ 
16   $global\_level := my\_level + 1$ 

```

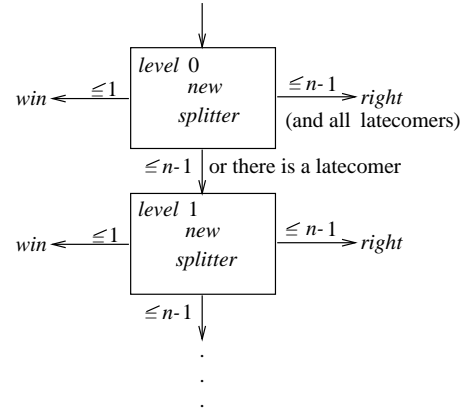


Figure 6: Adaptive deadlock-free mutual exclusion for unbounded concurrency.

of the next empty level, so that each process can join (or re-join) the competition starting from that level.

By the properties of the new splitter, when a process *wins* in a level, no process is active in any splitter with greater index. Every other active process will eventually move to the *right* in a splitter with equal or lesser index, and wait for the global pointer $global_level$ to be updated by the *winner*.

Measuring time complexity, exiting the critical section signals awaiting processes to proceed to a splitter at a new level. If n is the number of *early* processes at that splitter, then the next *winner* accesses at most $n + 1$ splitters before entering the critical section. This takes at most $O(n)$ time. Hence, the algorithm is adaptive.

Examining the details of the code, in the absence of contention, only seven accesses to shared atomic registers are needed before entering the critical section. ■

Process i 's program

Shared:

$global_level, cs_counter$: integer, initially both 0
 $x[0..\infty]$: array of integers (the initial values are immaterial)
 $b[0..\infty], y[0..\infty], z[0..\infty], try[1..\infty]$: array of boolean, initially all 0
 $winner_level$: integer (the initial value is immaterial)

Local:

my_level : integer (the initial value is immaterial)

```
0      my_level := global_level
0.1    try[i] := 1
1      start: x[my_level] := i
2      if y[my_level] then b[my_level] := 1 goto right fi
3      y[my_level] := 1
4      if x[my_level] ≠ i then await ((b[my_level] = 1) or (z[my_level] = 1))
5          if z[my_level] = 1 then goto right
6          else goto down fi
7      else z[my_level] := 1
8          if b[my_level] = 0 then goto win
9          else goto down fi fi
10     right: await ((my_level < global_level) or (try[i] = 0))
10.1   if try[i] = 0 then goto win fi          /* process i is helped */
11     my_level := global_level              /* no help, keep trying */
12     goto start
13     down: my_level := my_level + 1
14     goto start
15     win: critical section
15.1   if try[i] = 1 then winner_level := my_level fi    /* if i was not helped */
15.2   try[i] := 0                                     /* then post my_level */
15.3   cs_counter := cs_counter + 1                   /* help process Enum(cs_counter) */
15.4   if try[Enum(cs_counter)] = 1 then try[Enum(cs_counter)] := 0
16     else global_level := winner_level + 1 fi
```

Figure 7: Adaptive starvation-free mutual exclusion for unbounded concurrency.

3.2.3 Proof of Theorem 3.2

The main theorem of this section, Theorem 3.2, is an immediate consequence of the following.

Lemma 3.6 *The algorithm in Figure 7 is an adaptive, starvation-free mutual exclusion algorithm for unbounded concurrency.*

Proof: The adaptive deadlock-free algorithm of Figure 6 is easily modified using standard “helping” techniques to satisfy starvation-freedom. The details are in Figure 7. Process i sets a flag $try[i]$ when it leaves the remainder section, and before leaving the critical section, a process examines the flag of the next process (in some infinitely repeating enumeration $Enum(cs_counter)$ indexed by the number of times, $cs_counter$, that any process

Process i 's program

Shared:

$RaceOwner[1..\infty]$, $RaceOther[1..\infty]$, $Win[1..\infty]$, $Lose[1..\infty]$: arrays of boolean, initially all 0

Local:

$index$: integer, initially 1

```
1   $RaceOwner[i] := 1$ 
2  if  $RaceOther[i] = 0$  then  $Win[i] := 1$  else  $Lose[i] := 1$  fi
3  repeat forever
4       $RaceOther[index] := 1$ 
5      if  $RaceOwner[index] = 1$  then
6          await ( $Win[index] = 1$  or  $Lose[index] = 1$ )
7      fi
8      if  $Win[index] = 1$  then  $return(index)$  fi
9       $index := index + 1$ 
10 end repeat
```

Figure 8: (Non-adaptive) leader election for unbounded concurrency using bits

has entered the critical section) and grants the critical section if it determines the associated process is trying. Because the number of processes is infinite, an infinitely repeating enumeration is necessary instead of e.g. a round-robin schedule: a process helps others in the order given by an enumeration in which every process id appears infinitely often.

Starvation-freedom follows. Adaptivity follows just as in the deadlock-free algorithm: recall that in defining adaptivity for mutual exclusion algorithms, we measure the maximum time between releasing the critical section, until the critical section is re-entered. ■

Even simpler, standard modifications convert the deadlock-free mutual exclusion algorithm in Figure 6 to solve leader election or consensus: the first process to reach the critical section announces its id or the consensus value in an additional register. This completes the proof of Theorem 3.2. ■

3.2.4 A non-adaptive mutual exclusion algorithm using atomic bits

The adaptive mutual exclusion algorithm for unbounded concurrency we have just presented uses infinitely many infinite-sized registers. Infinite space is necessary—Alur and Taubenfeld showed that even without contention, $\Omega(\sqrt{n})$ shared bits must be accessed in any mutual exclusion algorithm for n processes [AT96]. Next we show that for a non-adaptive algorithm, it suffices to use only atomic bits.

Theorem 3.7 *There is a non-adaptive asymmetric solution to election, consensus and starvation-free mutual exclusion for unbounded concurrency using an infinite number of bits.*

Proof: Figure 8 presents a simple algorithm for election. Modification to solve consensus is trivial. As we explain, starvation-free mutual exclusion can be achieved using techniques similar to those in the previous algorithm. There are four bits associated with each process i , $RaceOwner[i]$, $RaceOther[i]$, $Win[i]$, and $Lose[i]$. The idea is that each process i uses the

two bits $RaceOwner[i]$ and $RaceOther[i]$ to race against the other processes. The minimum winner is the leader. The other two bits are used by the owner to signal whether it won or lost. The leader is the minimum i such that $Win[i]$ is set.

Termination: First, a simple argument shows each execution of the wait loop in line 6 eventually terminates. The first process i to perform the read in the second line will eventually set $Win[i]$. Every other process will either terminate before $index = i$, or see $RaceOwner[i]$ set, wait until $Win[i]$ is set, and terminate.

Agreement: In an infinite run, there is a minimum i such that $Win[i]$ is ever set. Every other process must see $RaceOwner[i]$ set, since process i reads $RaceOther[i] = 0$ in line 2 after setting $RaceOwner[i]$ in line 1. So every process will return i .

A fairly straightforward revision implements deadlock-free mutual exclusion. That is, use infinitely many copies of these binary data structures, one for every entrance to the critical section. Each copy x has a control bit, $Version[x]$. The leader in $Version[x]$ is the x 'th process to enter the critical section. The exit code is to set the control bit on the current version, signalling $Version[x + 1]$ is now current. To enter the critical section, a process finds the minimum copy with the control bit still open, contends there until it is the leader or (if it isn't) spins until the control bit is set on that copy, then tries again on the next.

To achieve starvation-free mutual exclusion, use the index to the current copy of the algorithm to determine who to help next in a (predetermined) infinitely repeating enumeration. See Figure 9 for details. (To implement the $LeaderElect(x)$ procedure, for $x \in \{1, \dots, \text{inf}\}$, process i runs the leader election algorithm of Figure 8 on copy x of the data structure. The $Enum(x)$ procedure returns the x 'th process id in an infinitely repeating enumeration.) ■

3.3 Relating participation and concurrency

Section 2 focused on the case in which the number of required participants, ℓ , is equal to the concurrency, c . Note that if $c < \ell$, the notion of required participation depends on termination of at least some of the first c participants. Indeed, the leader election algorithms of the previous section will work correctly for this case, as no process terminates until c have begun, and the leader is determined. If $\ell < c$, the situation is more complex. The following theorem unifies the results of Theorems 2.6 (participation is required) and 3.1 (participation is not required), demonstrating a relation between participation and concurrency in this case.

Theorem 3.8 *For concurrency c with $[\ell, \infty]$ -participation where $\ell < c$, $\log(c - \ell + 1) + 3$ registers are sufficient for solving election.*

Proof: The algorithm is a combination of the a single register assuming $[c, \infty]$ -participation (from Theorem 2.6), and of the one using $\log c + 1$ registers assuming $[1, \infty]$ -participation (from Theorem 3.1). First, we use the single register algorithm (Theorem 2.6) as a filter. Here, instead of choosing a single leader, up to $c - \ell + 1$ processes may be elected. These processes continue to the next level. To implement this, we slightly modify the single register algorithm. A process needs to wait only until it notices that ℓ (instead of c) processes are participating, and if it is the biggest among them it continues to the next level. Thus, at most $c - \ell + 1$ (and at least one) processes continue to the next level. In that level, they

Process i 's program

Shared:

$Version[1..\infty]$, $Try[1..\infty]$: arrays of boolean, initially all 0

Local:

v_index : integer, initially 1

$leader$: process id, initially nil

```
1  repeat forever
2      Remainder
3       $Try[i] := 1$ 
4      repeat forever
5          while  $Version[v\_index] = 1$  do  $v\_index := v\_index + 1$  od
6           $leader := LeaderElect(v\_index)$ 
7          if  $leader = i$  then break else
8              await  $((Version[v\_index] = 1) \vee (Try[i] = 0))$  fi
9              if  $Try[i] = 0$  then break fi
10     end repeat
11     Critical Section
12      $Try[i] := 0$ 
13     while  $Version[v\_index] = 1$  do  $v\_index := v\_index + 1$  od
14     if  $Try[Enum(v\_index)] = 1$  then  $Try[Enum(v\_index)] := 0$  fi
15         else  $Version[v\_index] := 1$  fi
16 end repeat
    /* help process  $Enum(v\_index)$  */, */
    /* or move to next version */
```

Figure 9: (Non-adaptive) starvation-free mutual exclusion for unbounded concurrency using bits.

compete using the algorithm of Theorem 3.1 (no change is needed in that algorithm), until one of them is elected. This level requires $\log(c - \ell + 1) + 1$ registers.

Finally, we use one more register, called *done*, for the leader to announce itself as the elected leader. Processes that are executing the single register algorithm (from Theorem 2.6) are also monitoring the *done* register and if they notice that a leader has been elected they terminate. ■

4 Test&set bits

A test&set bit is an object that may take the value 0 or 1, and is initially set to 1. It supports two operations: (1) a reset operation: write 0, and (2) a test&set operation: atomically assign 1 and return the old value. We first make the following observation:

Theorem 4.1 *An infinite number of atomic registers are necessary for implementing a test&set bit when concurrency is bounded, and sufficient when concurrency is unbounded.*

Proof: To prove the necessary condition assume to the contrary that such an implementation is possible with a fixed, finite number of atomic registers. This assumption, together with the observation that there exists a trivial deadlock-free mutual exclusion algorithm using a single test&set bit, implies that there is a solution to deadlock-free mutual exclusion with bounded concurrency using a finite number of atomic registers. This contradicts Theorem 3.1. The sufficient condition follows immediately from Theorem 3.7. ■

Theorem 4.2 *For solving starvation-free mutual exclusion, an infinite number of atomic bits and test&set bits are necessary and sufficient when the concurrency level is bounded.*

Proof: In [Pet94], it is proved that n atomic registers and test&set bits are necessary for solving the starvation-free mutual exclusion problem for n processes (with concurrency level n). This implies the necessary condition above. The sufficient condition follows immediately from Theorem 3.7. ■

4.1 Naming

The following theorem demonstrates that in certain cases a problem is solvable assuming bounded concurrency, but is not solvable assuming unbounded concurrency. The *wait-free naming* problem is to assign unique names to initially identical processes. After acquiring a name the process may later release it. A process terminates only after releasing the acquired name. A solution to the problem is required to be *wait-free*, that is, it should guarantee that every participating process will always be able to get a unique name and terminate in a finite number of steps regardless of the behavior of other processes (such as abnormal termination).

Theorem 4.3

1. *For bounded concurrency, an infinite number of T&S bits are necessary and sufficient for solving wait-free naming.*

2. For unbounded concurrency, there is no solution to wait-free naming, even when using an infinite number of T&S bits.

Proof: Part 1. The sufficient condition follows from the following simple algorithm, a variation of algorithms due to Alur and Taubenfeld [AT96] and to Moir and Anderson [MA95]. The algorithm uses an infinite number of bits with initial values 0, indexed 1,2,... Each process scans the bits, in order, starting with the first bit (i.e., bit number 1). At each step, the process applies the `test&set` operation, and either moves to the next bit if the returned value is 1, or stops when the returned value is 0. The process is assigned the name equal to the bit on which its (last) `test&set` operation returned 0. (Notice that wait-freedom is satisfied assuming bounded concurrency, however, it is not satisfied assuming unbounded concurrency since there is a run in which some process will always return 1's and never terminate.)

To prove the necessary condition, we use the following simple argument. Consider an algorithm which solves the problem. Pick a process, run it alone until it acquires a name. Then pick another process and run it alone until it acquires a name, and so on. Since the processes must be assigned unique names, also each state of the algorithm at the time a process terminates must be different from any other such state. Thus, the state space must be infinite and hence also the number of bits.

Part 2. Assume to the contrary that such an algorithm, say A , exists. Let p and p' be two processes. We reach a contradiction by constructing an infinite run α of A in which p and p' always execute the same steps, and hence can never acquire different names. We construct α inductively such that for each positive integer i , α has a prefix $\alpha(i)$ where:

1. p and p' have taken i steps each, and their steps are exactly the same.
2. none of the `test&set` operations of p and p' has ever returned 0.

We define $\alpha(0)$ to be the empty run. Assume that we have constructed $\alpha(i)$. We show how to extend it to $\alpha(i + 1)$. If the next operation of p and p' in $\alpha(i)$, is either a *reset* or a `test&set` on a bit which is already set to 1, then $\alpha(i + 1)$ is constructed by letting p and p' each take one step, and we are done. Otherwise, their next operation is a `test&set` operation on a bit, say b , which is set to 0. We first extend $\alpha(i)$ by letting some other processes take steps until b is set to 1 and only then let p and p' take one step each, as before.

Setting b to 1 using other processes is the tricky part. We use the following simple observation. If $i + 1$ identical processes apply the same operation to the same bit, then the returned values for at least i processes must be the same. Now consider the following extension of $\alpha(i)$. First we activate $i + 1$ new processes and let each one of them take one step. Since the processes are initially identical, each will apply the same operation to the same bit and hence at least i of them will get the same response. Thus, i processes are still identical. Next we activate these i identical processes and let each one take one step, and are guaranteed to end up with a set of at least $i - 1$ identical processes. We repeat this procedure i times, until one process, say q , has taken i steps. We claim that q is identical to both p and p' . This follows from the fact that, as for p and p' , process q has taken exactly i steps, and none of the `test&set` operations of q has ever returned 0.

We now activate q and let it set b to 1, and finally complete the construction of $\alpha(i + 1)$ by letting p and p' each take one step. ■

Process i 's program
Shared:
 $(ticket, valid)$: integers, initially $(0, 0)$
Local:
 $(ticket_i, valid_i)$: integers

```

1   $\langle (ticket_i, valid_i) := (ticket, valid) \rangle$ 
2   $ticket := ticket + 1 \rangle$ 
3  while  $ticket_i \neq valid_i$  do
4   $\langle valid_i := valid \rangle$  od
5  critical section
6   $\langle valid := valid + 1 \rangle$ 

```

Figure 10: FIFO mutual exclusion using read-modify-write: The ticket algorithm

5 Stronger Objects

5.1 RMW bits

A read-modify-write object supports a single operation, which atomically reads a value of a shared register, and based on the value read, computes some new value and assigns it back to the register. When assuming a fault-free model with required participation, many problems become solvable with small constant space. For example, as indicated in the second table in Section 1, it is trivial to implement consensus using a single, three-valued read-modify-write object. Mutual exclusion, however, is another matter:

Theorem 5.1

- *When only a finite number of RMW registers are used for solving starvation-free mutual exclusion with bounded concurrency, one of them must be of unbounded size.*
- *One unbounded size RMW register is sufficient for solving first-in, first-out adaptive symmetric mutual exclusion when the concurrency level is unbounded.*

Proof: It is shown in [B⁺82] that in a model which supports a read-modify-write operations any starvation-free mutex for n processes requires at least $\sqrt{2n} + \frac{1}{2}$ values. This result implies the necessary condition.

To prove the sufficient condition, we slightly modify the ticket algorithm from [F⁺89]. The algorithm in Figure 10 works as follows: A process wishing to enter its critical section takes the next available ticket and waits until its ticket becomes valid, at which point it can safely enter its critical section. When it exits, it discards its ticket and validates the next invalid ticket in order (even if this ticket has not been taken yet).

The shared register contains two fields, $ticket$ and $valid$, each a non-negative integer. The brackets $\langle \rangle$ are used to explicitly mark the beginning and end of exclusive access to the shared read-modify-write register. An execution of a bracketed section is considered as an atomic action. Each process first reads the value of the shared register, stores its components in local memory, and increments $ticket$ by one. At any later point, a process becomes the first in the “waiting line” if it learns (by inspecting $valid$) that its ticket number $ticket_i$ equals $valid$, in which case it can safely enter its critical section. ■

5.2 Semaphores

Given the results so far about starvation-free mutual exclusion, it is natural to ask whether it can be solved with a bounded number of semaphores. The answer, as presented in [FP87], is that using weak semaphores, it can be solved with small constant space for unbounded concurrency.

A binary semaphore is a shared objects that may take the values 0 or 1, and is initially set to 1. It support two operations, called P and V :

- When a process performs a P operation, if the value is 1, then the value is set to 0; otherwise, the process is blocked until the value is 1. Testing and decrementing the semaphore are executed atomically.
- When a process performs a V operation, the value is set to 1.

The result below refers to *weak* semaphores, in which a process that executes a V operation will not be the one to complete the next P operation on that semaphore, if another process has been blocked at that semaphore. Instead, one of the blocked processes is allowed to “pass” the semaphore to a blocked process.

Theorem 5.2 (Friedberg and Peterson [FP87]) *There is an adaptive symmetric solution to starvation-free mutual exclusion using two atomic bits and two weak semaphores, when the concurrency is unbounded.*

6 Discussion

We have explored how various assumptions about the number of processes, the concurrency level, and the number of participating processes, effect the design of shared memory algorithms. In particular, we have looked at the case of computing with infinitely many processes. There are many interesting questions that are left open: In the fault-free model using atomic registers only, are there interesting problems that can be solved for any finite and *a priori* known number of processes but cannot be solved for a finite and unknown number of processes (even with unbounded space)? Are there interesting problems that can be solved using a *finite* number of atomic registers by an infinite number of processes with bounded concurrency? Designing wait-free algorithms for problems such as collect, renaming and snapshot assuming unbounded concurrency are challenging problems. More recent work in this area demonstrates the existence of a hierarchy of computable tasks, determined by concurrency bounds [GMT2001].

Are there automatic transformations between some algorithms (such as symmetric algorithms) which work correctly for a finite number of processes, to algorithms which work correctly (perhaps assuming bounded concurrency) for infinitely many processes?

The relationship between algorithms for an infinite number of processes and adaptive algorithms is intriguing. Algorithms designed for an infinite number of processes are often adaptive, but are there natural conditions under which one type of algorithm implies the other? For example, consider the following symmetric algorithm for unbounded concurrency: Take a fixed enumeration, p_1, \dots of process id’s. Each process p_i writes the the id’s from the enumeration one at a time to a different shared register r_i , until process p_i ’s id is reached. This silly algorithm is not adaptive—the step complexity of even a solo execution is unbounded, as it depends on the position i of p_i ’s id in the enumeration.

Finally, are there algorithms that are correct assuming finite concurrency, but which fail if this constraint is not met? Can we, in general, modify such algorithms to guarantee at least safety (if not liveness) even when the constraints on the number of processes and the concurrency level are not met?

Acknowledgement We wish to thank an anonymous referee for many suggestions and corrections, and notably, catching an error in the conference version of the algorithm in Figure 3. We also thank the editor Faith Fich and Nancy Lynch for many suggestions and corrections.

References

- [AB2002] H. Attiya and V. Bortnikov. Adaptive and efficient mutual exclusion. *Distributed Computing*, 15(3):177–189, 2002.
- [AK2000] J.H. Anderson and Y. Kim. Adaptive mutual exclusion with local spinning. In *Proceedings of the 14th international symposium on distributed computing*, 2000.
- [AST99] Y. Afek, G. Stupp, and D. Touitou. Long-lived adaptive collect with applications. In *Proc. 40th IEEE Symp. on Foundations of Computer Science*, pages 262–272, October 1999.
- [AST2002] Y. Afek and G. Stupp and D. Touitou. Long-lived adaptive splitter with application. *Distributed Computing*, 15(2):67–86, 2002.
- [AT92] R. Alur and G. Taubenfeld. Results about fast mutual exclusion. In *Proceedings of the 13th IEEE Real-Time Systems Symposium*, pages 12–21, December 1992.
- [AT96] R. Alur and G. Taubenfeld. Contention-free complexity of shared memory algorithms. *Information and Computation* 126:1 (1996) 62–73. (Also in PODC 1994.)
- [AGM98] H. Attiya, A. Gorbach, and S. Moran. Computing in totally anonymous asynchronous shared memory systems. In *Proc. 12th International Symposium on Distributed Computing*, LNCS 1499:49-61, September 1998. Also in: *Information and Computation*, 173(2):162–18, March 2002.
- [ASS2002] J. Aspnes, G. Shah, and J. Shah. Wait-free consensus with infinite arrivals. In *Proc. 34th Annual Symp. on Theory of Computing*, 524–533, May 2002.
- [AY96] J-H. Yang and J.H. Anderson. Time/Contention Trade-Offs for Multiprocessor Synchronization. *Information and Computation*, 124(1):68–84, 1996.
- [BMT95] H. Brit, S. Moran, and G. Taubenfeld. Public data structures: counters as a special case. *Proc. 3rd Israel Symposium on Theory of Computing and Systems*, Tel Aviv, 98–110, January 1995. Also in: *Theoretical Computer Science*, 289(1):401–423, 2002.
- [B⁺82] J. E. Burns, P. Jackson, N. A. Lynch, M. J. Fischer, and G. L. Peterson. Data requirements for implementation of N -process mutual exclusion using a single shared variable. *Journal of the Association for Computing Machinery*, 29(1):183–205, 1982.

- [BL93] J. N. Burns and N. A. Lynch. Bounds on shared-memory for mutual exclusion. *Information and Computation*, 107(2):171–184, December 1993.
- [CM2002] G. Chocker and D. Malkhi. Active disk paxos with infinitely many processes. In *Proc. 21th ACM Symp. on Principles of Distributed Computing*, 78–87, July 2002.
- [CS94] M. Choy and A.K. Singh. Adaptive solutions to the mutual exclusion problem. *Distributed Computing*, 8(1):1–17, 1994.
- [Dij65] E. W. Dijkstra. Solution of a problem in concurrent programming control. *Communications of the ACM*, 8(9):569, 1965.
- [F⁺89] M. J. Fischer, N. A. Lynch, J. E. Burns, and A. Borodin. Distributed FIFO allocation of identical resources using small shared space. *ACM Trans. on Programming Languages and Systems*, 11(1):90–114, January 1989.
- [FLP85] M. J. Fischer, N. A. Lynch, and M. S. Paterson. Impossibility of distributed consensus with one faulty process. *Journal of the ACM*, 32(2):374–382, 1985.
- [FP87] S. A. Friedberg and G. L. Peterson. An efficient solution to the mutual exclusion problem using weak semaphores. *Information Processing Letters*, 25(5):343–347, 1987.
- [GK98] E. Gafni and E. Koutsoupias. On uniform protocols. <http://www.cs.ucla.edu/~eli/eli.html>, 1998.
- [Gafni2002] E. Gafni. A Simple Algorithmic Characterization of Uniform Solvability. In *43rd Symp. on Foundations of Computer Science*, 228–237, 2002.
- [GMT2001] E. Gafni, M. Merritt, and G. Taubenfeld. The concurrency hierarchy, and algorithms for unbounded concurrency. In *Proc. 20th ACM Symp. on Principles of Distributed Computing*, 161–169, August 2001.
- [Lam87] L. Lamport. A fast mutual exclusion algorithm. *ACM Trans. on Computer Systems*, 5(1):1–11, 1987.
- [LA87] M. C. Loui and H. Abu-Amara. Memory requirements for agreement among unreliable asynchronous processes. *Advances in Computing Research*, 4:163–183, 1987.
- [L⁺93] N. Lynch, M. Merritt, W. Weihl, and A. Fekete. *Atomic transactions*. 1993, Morgan Kaufmann.
- [Lyn96] N. Lynch. *Distributed algorithms*. 1996, Morgan Kaufmann.
- [MA95] M. Moir and J. Anderson. Wait-Free algorithms for fast, long-lived renaming, *Science of Computer Programming* 25(1):1–39, 1995.
- [MT93] M. Merritt and G. Taubenfeld. Speeding Lamport’s fast mutual exclusion algorithm. *Information Processing Letters*, 45:137–142, 1993. (Also published as an AT&T technical memorandum, May 1991.)

- [MT2000] M. Merritt and G. Taubenfeld. Computing with infinitely many processes. *Proceedings of the 14th International Symposium on Distributed Computing*, LNCS 1914, 164–178, October 2000.
- [MT2003] M. Merritt and G. Taubenfeld. Resilient Consensus for Infinitely Many Processes. *Proceedings of the 17th International Symposium on Distributed Computing*, LNCS 2648, 1–15, October 2003.
- [Pet94] G. L. Peterson. New bounds on mutual exclusion problems. Technical Report TR68, University of Rochester, February 1980 (Corrected, Nov. 1994).
- [PF77] G. L. Peterson and M. J. Fischer. Economical solutions for the critical section problem in a distributed system. In *Proc. 9th ACM Symp. on Theory of Computing*, pages 91–97, 1977.
- [SP89] E. Styer and G. L. Peterson. Tight bounds for shared memory symmetric mutual exclusion problems. In *Proc. 8th ACM Symp. on Principles of Distributed Computing*, 177–191, 1989.
- [Tau2004] G. Taubenfeld. The black-white bakery algorithm. In *18th international symposium on distributed computing*, LNCS 3274, 56–70, October 2004.