Lower Bounds in Distributed Computing



Petr Kuznetsov Telecom ParisTech (supported by Hagit Attiya, Technion)

> EPIT 2017 Île Porquerolles, France

© Richard Schneider

Lower and upper bounds

- If something can be computed, what are the costs?
- Complexity metrics:
 - ✓ time (e.g., number of steps)
 - ✓ space (number of shared memory locations)✓ typically asymptotic wrt n (system size)
- Lower bounds $\Omega(f(n))$ required cost
- Upper bounds O(f(n)) exhibited cost
- Tight lower bounds: f(n)=g(n)





This class

- Covering/valency arguments
 - ✓Time/space complexity of perturbable objects
 - ✓ Space complexity of obstruction-free consensus
- Information theory: encoder/decoder
 - ✓Total work complexity of mutual exclusion

A Tight Space Bound for Consensus

ment of Computer Science University of Toronto Canada

ABSTRACT

vising n-process randomized wait-free (and obstruction reprocess randomized wateries (and obstruction-nearsus protocols from registers all use at least n. In 1992, it was proved that such protocols must \overline{n}) registers. Recently, this was improved to $\Omega(n)$

CCS Concepts • Theory of computation \rightarrow Concurrence

Keywords Shared Memory Model Convensus Space Complexi

1. INTRODUCTION

Perhaps the most studied problem in the theory of dis tributed computing is the con usus problem, which require sees, each with an input value, to agree on a commo put value. An attractive application of the consensu oblem lies in implementing shared objects, such as stad queues. In particular, if there is a wdit-free pro where each pr ess decides in a finite numb its own steps, regardless of the speed or failure of oth n it is also possible to implement any shares nner Her91

ible to deterministcally solve w us shared me rs [LAA87] Ho n [AC08 AH90 AW96 CIL9/

STOC '16 June 18-21 2016 Combridge MA US

by the other many tight bounds were not into the moment space complexity of this problem. In 1992, Fich, H and Shavit proved a space lower bound of $\Omega(\sqrt{n})$ ters [FHS98]. All existing protocols use at least n re-[AH90, AW96]. Closing this gap has been a long open problem.

nymous (i.e. they have no identifiers) and mem they do not use local memory) [Zhu15]. At the s sing very interesting, different techniques, Gelas here are anonymous pro

here was even evidence suggesting the possibility of a pro-col using $O(\sqrt{n})$ space: Weak leader election is a closel

elated, but provably weaker, problem. rocesses must choose exactly one leader ow whether it has been vative protocol for weak leader election, using $O(\sqrt{n})$ regis ters, was obtained [GHHW13] a few years ago. Later, the ne authors imp wed this to $O(\log n)$ which is ont [GHHW15].

resolve the general case of the problem by proving that ters. Our lower b

An $\Omega(n \log n)$ Lower Bound on the Cost of Mutual Exclusion $A_{bstract}$



Part I

Covering



The idea



Processes need to communicate to get their work done Several processes may cover the same location

- An operation must write to enough distinct locations before terminating
- Otherwise, the operation is not visible

INFORMATION AND COMPUTATION 107, 171–184 (1993)

Bounds on Shared Memory for Mutual Exclusion*

JAMES E. BURNS

Georgia Institute of Technology, Atlanta, Georgia 30332

AND

NANCY A. LYNCH

Massachusetts Institute of Technology, Cambridge, Massachusetts 02139

The shared memory requirements of Dijkstra's mutual exclusion problem are examined. It is shown that n binary shared variables are necessary and sufficient to solve the problem of mutual exclusion with guaranteed global progress for n processes using only atomic reads and writes of shared variables for communication.



p₁,...,p_r

I.1: Perturbable complexity

SIAM J. COMPUT. Vol. 30, No. 2, pp. 438–456 © 2000 Society for Industrial and Applied Mathematics

TIME AND SPACE LOWER BOUNDS FOR NONBLOCKING IMPLEMENTATIONS*

PRASAD JAYANTI[†], KING TAN[†], AND SAM TOUEG[‡]

Abstract. We show the following time and space complexity lower bounds. Let \mathcal{I} be any randomized nonblocking *n*-process implementation of any object in set A from any combination of objects in set B, where $A = \{\text{increment}, \text{fetch}_{\text{add}}, \text{module } k \text{ counter} (\text{for any } k \geq n), \text{LL/SC bit}, k-\text{valued compare&swap (for any <math>k \geq n), \text{ single-writer snapshot}\}, \text{ and } B = \{\text{resettable consensus}\} \cup \{\text{historyless objects such as registers and swap registers}\}$. The space complexity of \mathcal{I} is at least n-1. Moreover, if \mathcal{I} is deterministic, both its time and space complexity are at least n-1. These lower bounds hold even if objects used in the implementation are of unbounded size.

This improves on some of the $\Omega(\sqrt{n})$ space complexity lower bounds of Fich, Herlihy, and Shavit [Proceedings of the 12th Annual ACM Symposium on Principles of Distributed Computing, Ithaca, NY, 1993, pp. 241–249; J. Assoc. Comput. Mach., 45 (1998), pp. 843–862]. It also shows the near optimality of some known wait-free implementations in terms of space complexity.

Key words. asynchronous shared memory algorithms, nonblocking, wait-free, synchronization, randomized shared object implementations, space complexity, time complexity, lower bounds

AMS subject classifications. 68Q17, 68W15

PII. S0097539797317299

1. Introduction. Nonblocking and wait-free implementations of shared objects have been the subject of much research. While there have been several results on when such implementations are feasible and when they are not, results establishing their intrinsic time and space requirements are relatively scarce, especially for randomized implementations. In this paper, we present a technique by which one can obtain a linear lower bound on the space complexity of several randomized nonblocking implementations. The technique also yields a linear lower bound on the time complexity of several deterministic nonblocking implementations.

Specifically, our results are as follows. Let \mathcal{I} be any randomized nonblocking n-process implementation of any object in set A from any combination of objects in set B, where $A = \{$ increment, fetch&add, modulo k counter (for any $k \geq 2n$), LL/SC bit, k-valued compare&swap (for any $k \geq n$), single-writer snapshot}, and $B = \{$ resettable consensus $\} \cup \{$ historyless objects $\}$. (Roughly speaking, an object is historyless if each of its operations either does not affect the state of the object or overwrites the previously applied operations. Examples include registers, test and set objects, and swap registers.) The space complexity of \mathcal{I} is at least n - 1. Moreover, if \mathcal{I} is deterministic, both its time and space complexity are at least n - 1. These lower bounds hold even if objects used in the implementation are of unbounded size.

Some of the results in this paper improve known lower bounds, while others are completely new. In particular, Fich, Herlihy, and Shavit proved a $\Omega(\sqrt{n})$ space

The result (simplified)

Any linearizable obstruction-freedom implementation of a counter from readwrite registers has

- \geq n-1 space complexity
- \geq n-1 solo step complexity

Space and (solo-step) time complexity is $\Theta(n)$ [Attiya et al., JACM'09]



Obstruction-freedom

An operation is guaranteed to return if it runs in isolation (no step contention) for sufficiently long



OF read-write implementations: consensus, CAS, counters,...

Model assumptions

- Processes *II*={p₁,...,p_n} communicate via reading and writing to shared base objects
- Every process is assigned a deterministic (counter) algorithm
- Every process runs inc() operations, one after another
 - ✓ The state of the system is determined by a schedule --- a sequence in Π^*



k=n-1 - we are done!

Base case

- k=0
- $\alpha_0 = \beta_0 = \gamma_k = \epsilon$ (empty schedules)

Suppose the hypothesis holds for $0 \le k \le n-1$

Induction step: λ_k can "perturb" p_n



Claim: p_n must access some B_{k+1} outside $\{B_1, \dots, B_k\}$ before returning in γ

Suppose not: p_n returns v only accessing $B_1, ..., B_k$

 p_n can be "perturbed" by λ_k



- α_{k+1} extend α_k until p_{k+1} is about to write to some B_{k+1} outside $\{B_1, ..., B_k\}$
- β_k=p₁,...,p_{k+1}
- γ_{k+1} extend γ_k until p_n accesses B_{k+1}

Perturbable objects

There exists an assignment of operations such that for every schedule $\alpha\beta\gamma$ such that

- α and β do not contain p_n
- β is by a proper subset of {p₁...p_{n-1}}
- $\gamma \in p_n^*$ and p_n runs exactly one operation in $\alpha\beta\gamma$
- For some p_I∈{p₁,...,p_{n-1}}-PSET(β), ∃ λ ∈ p_I*: p_n does not complete its operation or returns a different response in αλβα and αβα

Ω(n) time/space hold for *randomized* implementations of *perturbable* objects (CAS, counters, atomic snapshots)
 from historyless primitives (read-write, swap, ...)

I.2: Space complexity of consensus: valence and covering

A Tight Space Bound for Consensus

Leqi Zhu Department of Computer Science University of Toronto Canada lezhu@cs.toronto.edu

ABSTRACT

Existing n-process randomized wait-free (and obstructionfree) consensus protocols from registers all use at least n registers. In 1992, it was proved that such protocols must use $\Omega(\sqrt{n})$ registers. Recently, this was improved to $\Omega(n)$ registers in the anonymous setting, where processes do not have identifiers. Closing the gap in the general case, however, remained an open problem. We resolve this problem by proving that every randomized wait-free (or obstructionfree) consensus protocol for n processes must use at least n - 1 registers.

CCS Concepts

Keywords

\bullet Theory of computation \rightarrow Concurrency

Shared Memory Model, Consensus, Space Complexity

1. INTRODUCTION

Perhaps the most studied problem in the theory of distributed computing is the consensus problem, which requires *n* processes, each with an input value, to agree on a common output value. An attractive application of the consensus problem lies in implementing shared objects, such as stacks or queues. In particular, if there is a *wait-free* protocol for consensus, where each process decides in a finite number of its own steps, regardless of the speed or failure of other processes, then it is also possible to implement any shared object in a wait-free manner [Her91].

It is impossible to deterministcally solve wait-free consensus in an asynchronous shared memory system, where processes communicate by reading and writing shared memory locations, called *registers* [LAA87]. However, it is possible using randomization [AC08, AH90, AW96, CLB4]. Asymptotically tight bounds are known for the total number of steps taken by all processes [AC08]. On the other hand, tight bounds were not known for the space complexity of this problem. In 1992, Fich, Herliky, and Shavit proved a space lower bound of $\Omega(\sqrt{n})$ registers [FHS98]. All existing protocols use at least *n* registers [AH90, AW96]. Closing this gap has been a longstanding open problem.

Recently, we proved matching upper and lower bounds of n registers for a restricted class of protocols, where processes are anonymous (i.e. they have no identifiers) and memoryless (i.e. they do not use local memory) [Zhu15]. At the same time, using very interesting, different techniques, Gelashvili proved a lower bound of $\Omega(n)$ registers for protocols with anonymous processes, without the memoryless assumption [Gel15]. Since there are anonymous protocols that use nregisters [BRS15, Zhu15], the bound is tight. Thus, the anonymous case of the problem is resolved to within a constant factor.

The general case of the problem, however, remained open. There was even evidence suggesting the possibility of a protocol using $O(\sqrt{n})$ space: Weak leader election is a closely related, but provably weaker, problem. In this problem, processes must choose exactly one leader, but each process only needs to know whether it has been chosen. An innovative protocol for weak leader election, using $O(\sqrt{n})$ registers, was obtained [GHHW13] a few years ago. Later, the same authors improved this to $O(\log n)$, which is optimal [GHHW15].

Our contribution.

We resolve the general case of the problem by proving that any consensus protocol for n processes in an asynchronous system uses at least n - 1 registers. Our lower bound uses a more refined notion of valency (introduced in [FLP85]) combined with a covering argument (introduced in [BL93]). As in [FH398, Gel15, Zhu15], the bound holds even if the registers are of unbounded size.

The lower bound shows that consensus is, fundamentally, a communication problem. In particular, having large registers cannot compensate for having too few registers. Since

The result

Space complexity of any obstructionfree binary consensus implementation is n-1

Θ(n) space complexity [AGHK09,BRS15] (applies to randomized consensus too)





Refined model assumptions

- Processes *I*={p₁,...,p_n} communicate via reading and writing to shared base objects
- Every process is assigned a deterministic algorithm
- Configuration
 - ✓ Local state for each process
 - ✓ State of each register
- Initial configuration
 - ✓ Input assignment
 - ✓ Registers in initial states
- Configuration C and schedule $\alpha \in \Pi^*$ define a configuration $\mathbf{C}\alpha$

Valence of a configuration

C is *v-valent* (for v in {0,1}) if v is decided in every extension of C

C is *bivalent* if both 0 and 1 can be decided in extensions of C

- Every configuration is 0-valent, or 1-valent, or bivalent.
- If some process decides v in C, then C is v-valent
- No process can decide in a bivalent run

Refined valence of a configuration

Let C be a configuration and $P \subseteq \Pi$

P is *v*-valent from C (for v in $\{0,1\}$) if v is decided in every P-only extension C α

P is *bivalent from C* if both 0 and 1 can be decided in P-only extensions of C

For all C and non-empty P:

- P is 0-valent, or 1-valent, or bivalent from C
- If P is v-valent in C, then any non-empty P'⊆ P is v-valent from C

Valence: more properties

Lemma 1 If P (IPI \geq 3) is bivalent from C, then there exists a Ponly schedule ϕ and $z \in P$ such that P-{z} is bivalent from C ϕ

Proof: Take z_1 , $z_2 \in P$, let $Q_1 = P - \{z_1\}$ and $Q_2 = P - \{z_2\}$, $Q_1 \cap Q_2 \neq \emptyset$



Valence and read-write: more properties

- **Lemma 2** Let P be bivalent from C, and β be a block write (in C) by R \subseteq P. Then every deciding schedule from C by z not in P should contain a write not covered by R in C.
- **Lemma 3** Let P-R be bivalent from C, and β be a block write by non-empty R \subseteq P. If Q=P-R is bivalent from C, then there exists Q-only schedule φ and $q \in Q$ such that $R \cup \{q\}$ is bivalent from $C\varphi\beta$.



The goal

If P, IPI \geq 2, is bivalent from C, then there exists P-only α and Q \subseteq P, |Q|=2, such that

- Q is bivalent from $C\alpha$
- Every process in P-Q covers a distinct register in $\mathbf{C}\alpha$

Proof by induction on IPI:

 Base case IPI=2: any initial configuration from which some P={p,q} is bivalent



For IPI=n we are done: n-2 distinct registers are covered

Induction step

P (of size \geq 3) is bivalent from C, the claim hold for subsets of size IPI-1

By Lemma 1, for some $z \in P$ and P-only γ , P-{z} is bivalent from D=C γ

By induction hypothesis, there exist D_0 , extension of D, and a pair $Q_0 \subseteq P$ -{z} such that

- Q_0 is bivalent from D_0
- P-{z}-Q₀ cover distinct registers in D₀



Induction step (ctd.)

Applying Lemma 3 and induction hypothesis repeatedly we get an infinite (P-{z})-only extension $D_0 \rightarrow D_1 \rightarrow D_3 \rightarrow ...$

- D_i: some pair $Q \subseteq P\$ is bivalent and the rest cover IPI-2 registers



Induction step (ctd.)

But there are finitely many registers! Some D_i and D_i cover the same set of IPI-3!



Q_i is bivalent P-{z}-Q_i cover Q_{i+1} is bivalent P-{z}-Q_{i+1} cover the same set

Induction step (ctd.)

Run z from $D_i \phi_i$ until it decides: $D_i \phi_i \zeta$

z must write to a not covered register (Lemma 2), stop just before: $\mathsf{D}_i\phi_i\zeta$

No (P-{z}) extension of $D_i \phi_i \zeta' \beta_i$ can see the difference from $D_i \phi_i \beta_i$



Continue by P-{z} until D'_{j} and get a cover by IPI-2 processes! For k=n, n-2 distinct registers are covered!

Finally!

Take k=n Let {p,q} be bivalent in D_j' {p,q}-only extension of D_i' must write to a non-covered register



n-1 registers covered!

Wrapping up: covering/valence

- Perturbable objects (CAS, counters, AS,...)
 ✓ "pure covering" assuming long-lived (perturbable) operations: time&space
- One-shot (non-perturbable) consensus
 ✓Covering&valence: space
- The proofs are about constructing a (worstcase) run

And now for something completely different



Part II

Information theory



$\Omega(n \log n)$ total work in mutual exclusion: The encoder/decoder argument

An $\Omega(n \log n)$ Lower Bound on the Cost of Mutual Exclusion

Rui Fan MIT CSAIL rfan@theory.csail.mit.edu

Nancy Lynch MIT CSAIL lynch@theory.csail.mit.edu

Abstract

non-busywaiting memory accesses by any deterministic algorithm solving n process mutual exclusion that communicates via shared registers. The cost of the algorithm is measured in the *state change* cost model, a variation of the cache coherent model. Our bound is tight in this model. We introduce a novel information theoretic proof technique. We first establish a lower bound on the information needed by processes to solve mutual exclusion. Then we relate the amount our proof technique is flexible and intuitive, and may be applied to a variety of other problems and system models.

1 Introduction

rithm is charged only for performing shared memory operations causing a process to change its state. Let We prove an $\Omega(n \log n)$ lower bound on the number of a *canonical execution* consist of n different processes, each of which enters the critical section exactly once. We prove that any deterministic mutex algorithm using registers must incur a cost of $\Omega(n \log n)$ in some canonical execution. This lower bound is tight, as the algorithm of Yang and Anderson [13] has $O(n \log n)$ cost in all canonical executions with our cost measure. To prove the result, we introduce a novel technique which is *information theoretic* in nature. We first argue that in each canonical execution, processes of information processes can acquire through shared need to cumulatively acquire a certain amount of inmemory accesses to the cost they incur. We believe formation. We then relate the amount of information processes can obtain by accessing shared memory to the cost of those accesses, to obtain a lower bound on the cost of the mutex algorithm. Our technique can be extended to show the same lower bound when processes are allowed access to comparison-based shared memory objects, in addition to registers. Further-

The result

Total work of any n-process mutual exclusion algorithm is

 $\Omega(n\log n)$

The number of (non busy-waiting) memory accesses performed by $p_1,..,p_n$ to enter CS

- $\checkmark \approx$ remote memory references (RMRs) in CC and DSM memory models [GW12,...]
- Tight [Yang&Anderson, 95]:

 $\Theta(n\log n)$



Holds even for stronger primitives (CAS, …)

Mutual exclusion

- No two processes are in their critical sections (CS) at the same time
- Deadlock-freedom: at least one process in its trying section (TS) eventually enters its CS
 - ✓ assuming no process fails or stays in its CS forever



Peterson's mutual exclusion

```
// code for process i that wishes to enter CS
for (m = 0; m < n-1; m++) {
    level[i] = m;
    waiting[m] = i;
    while(waiting[m] == i &&(exists k ≠ i: level[k] ≥ m)) {
        // busy wait
    }
}
// critical section
level[i] = -1; // exit section
Total work O(n<sup>3</sup>)
```

The idea: acquiring information incurs costs

- A canonical execution: every process enters CS exactly once
- Processes (cumulatively) must learn about the order of CSs
- Getting O(C) bits of information \equiv performing O(C) work

n! distinct orders $\equiv \Omega(n \log n)$ work

Visibility graph of a canonical run

 p_i "sees" $p_j \equiv$ the CS of p_i is causally preceded by CS of p_i

p_i left CS before p_i started its CS

Claim: in a canonical run, for every p_i and $\mathsf{p}_j,$ at least one sees the other



Proof outline

- 1. **Construction step**: for each permutation $\pi = \pi_1, \dots, \pi_n$, build a (canonical) run α_{π} with order π of CS accesses
- 2. Encoding step: for each α_{π} , produce a binary string E_{π} of length $O(cost(\alpha_{\pi}))$
- 3. Decoding step: reproduce α_{π} given E_{π}



- {E_{π}} is a code of { α_{π} }
- Some codeword E_{π}^{n} has length $\Omega(n \log n)$
- The cost of π is $\Omega(n \log n)$

1. Construction



 α_{π} is constructed iteratively:

- α_1 : p_{π_1} enters CS and exits
- From α_i to α_{i+1}: add a complete run of p_{πi+1} so that no process, so that p_{π1},...,p_{πi} do not see it ✓ The trickiest part: maintain a partial order on metasteps (M_i, ≤_i)
 Metastep (on the same register):

Set of writes

writes

Winning write

Signature: counts of reads and

• (M_n,
$$\preceq_n$$
) $\Rightarrow \alpha_\pi$

2. Encoding





$\begin{array}{l} (\mathsf{M}_n, \preceq_n) \to \mathsf{E}_{\pi} \\ \text{Only metasteps are encoded} \\ \text{The cost of a metastep with k processes is O(k)} \\ \mathsf{E}_{\pi} \text{ uses O(k) bits per metastep with k processes} \end{array}$





Iteratively compute the minimal "unexecuted" step:

- Given the algorithm and the current state of each process
- Compute the metastep's composition

 \Rightarrow (M_n, \leq_n)

Wrapping up

- Not everything (computable sequentially) can be computed in a distributed way
 - ✓ Computable problems are subject to complexity bounds
- There are many lower bounds but fewer techniques
 - ✓Covering/valence/potential functions
 - ✓ Information theory
 - ✓Combinatorial arguments



(Some) open questions

- From n-1 to n for consensus/perturbable objects
- (Solo) time complexity for OF consensus
- Time/space for k-set consensus
 - \checkmark Only anonymous case is explored
 - ✓ $\Omega(\operatorname{N})$ for m-OF k-set consensus
 - $\checkmark \Theta(n+m-k)$ for repeated m-OF k-set consensus
- Other system parameters
 - ✓ Progress guarantees
 - ✓ Contention
- Other metrics/beyond worst-case
- Computability bounds
 - Tight impossibility results: what is the weakest model for a given pb?

Merci beaucoup! Questions?



Induction hypothesis

For all k=0,...,n-1, there exist schedules $\alpha_{\rm k},\,\beta_{\rm k},\gamma_{\rm k}$ such that:

- α_k and β_k do not contain p_n
- For k=0, $\beta_k = \epsilon$ (empty), for k>0, $\beta_k = p_1, \dots, p_k$
- $\gamma_{\mathsf{k}} \in \mathsf{p_n}^*$
- The set of objects accessed by ${\rm p_n}$ in $\alpha_{\rm k}\beta_{\rm k}\gamma_{\rm k}$ is $\{{\rm B_1,...,B_k}\}$
- p_n performs at most one inc() in $\alpha_k \beta_k \gamma_k$
- Let λ_k be any schedule of $\{p_{k+1},...,p_{n-1}\}$. Then p_n cannot distinguish $\alpha_k \lambda_k \beta_k \alpha_k$ and $\alpha_k \beta_k \alpha_k$

k=n-1 - we are done!

Induction step

- α_k and β_k do not contain p_n
- β_k is a block write by $p_1,...,p_k$ to registers $B_1,...,B_k$
- $\gamma_{\mathbf{k}} \in \mathbf{p_n}^*$ and $\mathbf{p_n}$ performs at most one operation in $\alpha_{\mathbf{k}}\beta_{\mathbf{k}}\gamma_{\mathbf{k}}$
- Let λ_k be any schedule of $\{p_{k+1},...,p_{n-1}\}$. Then p_n cannot distinguish $\alpha_k \lambda_k \beta_k \alpha_k$ and $\alpha_k \beta_k \alpha_k$

Proof of Lemma 1

- Take $z_1, z_2 \in P$
- Let $Q_1 = P \{z_1\}$ and $Q_2 = P \{z_2\}$, $Q_1 \cap Q_2 \neq \emptyset$
- Let $Q_1 \cap Q_2$ decides $v \in \{0,1\}$ from C

 \checkmark If some Q_i decides 1-v from C - we are done

C is bivalent: 1-v is decided in C ψ

There must be a critical step in ψ : from C ψ ' ("only v is decided by Q₁ and Q₂") to C ψ ' δ ("some Q_i decides 1-v")

Let z (taking this step δ) be in Q₁:

- Q_1 is v-valent from $C\psi$ 'z
- Q_2 is bivalent from $C\psi'z$