

How to Elect a Leader Faster than a Tournament

Dan Alistarh
Microsoft Research

Rati Gelashvili
MIT

Adrian Vladu
MIT

Abstract

The problem of electing a leader from among n contenders is one of the fundamental questions in distributed computing. In its simplest formulation, the task is as follows: given n processors, all participants must eventually return a *win* or *lose* indication, such that a single contender may *win*. Despite a considerable amount of work on leader election, the following question is still open: can we elect a leader in an asynchronous fault-prone system faster than just running a $\Theta(\log n)$ -time tournament, against the classic omniscient strong adversary?

In this paper, we answer this question in the affirmative, improving on a decades-old upper bound. We introduce two new algorithmic ideas to reduce the time complexity of electing a leader to $O(\log^* n)$ rounds of communication, using $O(n^2)$ point-to-point messages. A non-trivial application of our algorithm is a new upper bound for the *tight renaming* problem, assigning n items to the n participants in expected $O(\log^2 n)$ time and $O(n^2)$ messages. We complement our results with lower bound of $\Omega(n^2)$ messages for solving these two problems, closing the question of their message complexity.

1 Introduction

The problem of picking a leader from among a set of n contenders in a fault-prone system is among the most well-studied questions in distributed computing. In its simplest form, *leader election (test-and-set)* [AGTV92] is stated as follows. Given n participating processors, each of the contenders must eventually return either a *win* or *lose* indication, with the property that a single participant may *win*. Leader election is one of a few canonical problems, or *tasks*, whose solvability and complexity are the focus of distributed computing theory, along with *consensus (agreement)* [LSP82, PSL80], *mutual exclusion* [Dij65], *renaming* [ABND⁺90], or *task allocation (do-all)* [KS92]. These problems are usually considered in *asynchronous* models, such as message-passing or shared-memory [Lyn97].

We focus on leader election in the *asynchronous message-passing* model, in which each of n processors is connected to every other processor via a point-to-point channel. Communication is asynchronous, i.e., messages can be arbitrarily delayed. The system is controlled by a *strong (adaptive) adversary*, which can examine local state, including random coin flips, and crash $t < n/2$ of the participants at any point during the computation. The natural complexity metrics are *message complexity*, i.e., total number of messages sent by the protocol, and *round (time) complexity*, i.e., number of send-receive steps executed by a processor before returning. A key difficulty is that rounds of communication are not synchronized across processors.

Many fundamental results in distributed computing are related to the complexity of canonical tasks in asynchronous models. For example, Fisher, Lynch, and Patterson [FLP85] showed that it is impossible to solve consensus deterministically in an asynchronous system if one of the n participants may fail by crashing. This deterministic impossibility extends to leader election [Her91]. Since the publication of the FLP result, a tremendous amount of research effort has been invested into overcoming this impossibility for canonical tasks. Seminal work by Ben-Or [BO83] showed that relaxing the problem specification to allow probabilistic termination can be a very used to circumvent FLP, and to obtain efficient distributed algorithms.

Consequently, the past three decades have seen a continuous quest to improve the randomized upper and lower bounds for canonical tasks, and in fact tight (or almost tight) complexity bounds are now known for consensus [AC08, AAKS14], mutual exclusion [HW09, HW10, GW12b], renaming [AACH⁺13], and task allocation [BKRS96, ABGG12] against a strong adversary.

For leader election, the situation is different. The fastest known solution is more than two decades old [AGTV92], and is a *tournament tree*: pair up the participants into two-processor “matches,” decided by two-processor randomized consensus; winners continue to compete, while losers drop out, until a single winner prevails. Time complexity is logarithmic, as the winner has to communicate at each tree level. No time lower bounds are known. Despite significant recent interest [AAG⁺10, AA11, GW12a], the question of whether a tournament is optimal when elect a leader against a strong adversary is surprisingly still open.

Contribution. In this paper, we show that it is possible to break the logarithmic barrier in the classic asynchronous message-passing model, against an adaptive adversary. We present a new randomized algorithm which elects a leader in expected $O(\log^* n)$ time, sending $O(n^2)$ messages.

The algorithm is based on two new ideas, which we briefly describe below. The general structure is rather simple: computation occurs in *phases*, where each phase is designed to drop as many participants as possible, while ensuring that at least one processor survives. Consider a simple implementation: each processor flips a biased coin at the beginning of the phase, to decide whether to give up (value 0) or continue (value 1), and communicates its choice to others. If at least one processor out of the n_r participants in phase r flips 1, all processors which flipped 0 can safely drop from contention. We could aim for $o(\log n)$ iterations by setting the probabilities to obtain less than a constant fraction of survivors in each phase. Unfortunately, a strong adversary easily breaks such a strategy: since it can see the flips, it can schedule all the processors that flipped 0 to complete the phase *before* any processor that flipped 1, forcing *everyone* to continue.

Techniques. Our first algorithmic idea is a way to *hide* the processor coin flips during the phase, handicapping the adaptive adversary. In each phase, each processor first takes a “poison pill” (moves to *commit*

state), and broadcasts this to all other processors. The processor then flips a biased local coin to decide whether to drop out of contention (*low* priority) or to take an “antidote” (*high* priority), broadcasts its new state, and checks the states of other processors. Crucially, if it has flipped low priority, and sees *any other processor* either in *commit* state or in *high priority* state, the processor returns *lose*. Otherwise, it survives to the next phase.

The above mechanics guarantee at least one survivor, (in an unlikely event where all processors flip *low* priority, they all survive.) but can lead to few survivors in each phase. The insight is that, to ensure many survivors, the adversary must examine the processors’ coin flips. But to do so, the adversary must first allow it to take the poison pill (state *commit*). Crucially, any low-priority processor observing this *commit* state automatically drops out. We prove that, because of this catch-22, the adversarial scheduler can do no more than to let processors execute each phase sequentially, one-by-one, hoping that the first processor flipping high priority, which eliminates all later low-priority participants, comes as late as possible in the sequence. Now we can bias the flips such that a group of at most $O(\sqrt{n_r})$ processors survive because they flipped high priority, and $O(\sqrt{n_r})$ processors survive because they did not observe any high priority. This choice of bias seems hard to improve, as it yields the perfect balance between the sizes of the two groups of survivors.

Our second algorithmic idea breaks this roadblock. Consider two extreme scenarios for a phase: first when all participants communicate with each other, leading to similar views and second, when processors see fragmented views, observing just a subset of other processors. In the first case, each processor can safely set a low probability of surviving. This does not work in the second case since processor views have a lot of variance, but this variance could be used to break symmetry. Our technical argument combines these two strategies such that we obtain at most $O(\log^2 n_r)$ expected survivors in a phase, under *any* scheduling.

The final algorithm has additional useful properties. It is *adaptive*, meaning that, if $k \leq n$ processors participate, its complexity becomes $O(\log^* k)$. Moreover, since most participants drop in the first round of broadcast, the message complexity is $O(kn)$, which we prove is asymptotically optimal.

Renaming. We build on these properties to design a message-optimal algorithm for *strong renaming*, which assigns distinct items (or names) labeled from 1 to n to the n processors, using expected $O(n^2)$ messages and $O(\log^2 n)$ communication rounds. We employ a simple strategy: each processor repeatedly picks a random name that it sees as available, announces it, and competes for it via an instance of leader election. If the processor wins, it returns the name; otherwise, it tries again. The algorithm can be seen as a balls-into-bins game, in which n balls are the processors and bins are the names. We need to characterize two parameters: the maximum number of trials by a single processor, and the maximum contention on a single bin, as they are linked with message and round complexity. The critical difficulty is that, since rounds are not synchronised, the bin occupancy views perceived by the processors are effectively under adversarial control and out-of-date or incoherent views can lead to wasted trials and increased contention on the bins.

Our task is to prove that, in fact, this balls-into-bins process is robust to the correlations and skews in the trial probability distributions caused by asynchrony. Our approach is to carefully bound the evolution of processors’ views and their trial distributions as more and more trials are performed. Roughly, for $j \geq 1$, we split the execution into time intervals, where at most $n/2^{j-1}$ names are available at the beginning of the interval j , and focus on bounding the number of wasted trials in each interval. The main technical difficulty that we overcome is that the views corresponding to these trials could be highly correlated, as the adversary may delay messages to increase the probability of a collision.

Lower bound. We match the message complexity of our algorithms with an $\Omega(n^2)$ lower bound on the expected message complexity of any leader election or renaming algorithm. The intuition behind the bound is that no processor should be able to decide without receiving a message, as the others might have already elected a winner; since the adversary can fail up to $n/2$ processors, it should be able to force each processor to either send or receive $n/2$ messages. However, this intuition is not entirely correct, as groups of processors could employ complex gossip-like message distribution strategies to guarantee that at least *some* processors

receive *some* messages while keeping the total message count $o(n^2)$. We thwart such strategies with a non-trivial indistinguishability argument, showing that in fact there must exist a group of $\Theta(n)$ processors, each of which either sends or receives a total of $\Theta(n)$ messages. A similar argument yields the $\Omega(n^2)$ lower bound for renaming, and in fact for any object with strongly non-commutative operations [AGH⁺11].

In sum, we reduce the time complexity of electing a leader from among k participants against an adaptive adversary from logarithmic to $O(\log^* k)$, and consequently give the first tight bounds on the message complexity of leader election and renaming in asynchronous message passing. The low round complexity of our leader election algorithm has implications for the power of topological lower bound techniques, e.g. [HS99], against randomized algorithms, which we discuss in Section 6.

Related Work. We focus on previous work on the complexity of randomized leader election¹ and renaming, most of which considered the asynchronous shared-memory. However, one option is to emulate efficient shared-memory solutions via simulations between shared-memory and message-passing [ABND95]. This preserves time complexity, but communication may be increased by at most a linear factor.

We classify previous solutions according to their adversarial model. Against a strong adversary, the fastest known leader election algorithm is the tournament tree of Afek et al. [AGTV92], whose contention-adaptive variant was given in [AAG⁺10]. For n participants, these algorithms require $\Theta(\log n)$ communication rounds, and $\Theta(n^2 \log n)$ messages using a careful simulation. *PoisonPill* is contention-adaptive, improves time complexity (more than) exponentially, and gives tight message complexity bounds.

For renaming, the fastest known shared-memory algorithm [AACH⁺13] can be simulated using $O(\log n)$ communication rounds, and $\Theta(n^2 \log n)$ messages. (The latter bounds are obtained by simulating an AKS sorting network [AKS83]; constructible solutions pay an extra logarithmic factor in both measures.) Our balls-into-bins approach trades simplicity and optimal message complexity for an extra logarithmic factor in the time complexity. Reference [AAG⁺10] uses a similar, but much simpler balls-into-bins approach for renaming, where each processor tries all the names, in random order, until acquiring some name. Despite the similarity, this algorithm has expected time complexity $\Omega(n)$, as a late processor may try out a linear number of spots before succeeding.

References [AA11, GW12a] considered the complexity of leader election against a weak (oblivious) adversary, which fixes the schedule in advance. The structure of splitting the computation into sifting rounds, eliminating more than a constant factor of the participants per round, was introduced in [AA11], where the authors give an algorithm with $O(\log \log n)$ time complexity. Giakkoupis and Woelfel [GW12a] improved this to $O(\log^* k)$ rounds, where k is the number of participants. These algorithms yield the same bounds in asynchronous message-passing, however, their complexity bounds only hold against a weak adversary.

The *consensus* problem can be stated as leader election if we ask processors to return the *identifier* of the winner, as opposed to a win/lose indication. As such, consensus solves leader election, but not vice-versa [Her91]. In fact, randomized consensus has $\Omega(n)$ round complexity [AC08]. Recent work [AAKS14] considered the message complexity of randomized consensus in the same model, achieving $O(n^2 \log^2 n)$ message complexity, and $O(n \log^3 n)$ round complexity, using fundamentally different techniques.

2 Definitions and Notation

System Model. We consider the classic asynchronous message-passing model [ABND95]. Here, n processors communicate with each other by sending *messages* through *channels*. There is one channel from each processor to every other processor; the channel from i to j is independent from the channel from j to i . Messages can be arbitrarily delayed by a channel, but do not get corrupted.

Computations are modeled as sequences of steps of the processors. At each step of a processor, first, some messages from incoming channels may be delivered to it. Then, unless the processor is *faulty*, it

¹Some older references, e.g. [AGTV92], employ the name *test-and-set* exclusively for this task, and use leader election for the consensus (agreement) problem, while more recent ones [GW12a] equate test-and-set and leader election.

can perform local computations and send new messages. More formally, a processor is *non-faulty*, if it is allowed steps to perform local computations and send messages infinitely often and if all messages it sends are eventually delivered. Every message delivery corresponds to a step. For convenience, we assume that messages are also delivered to *faulty* processors; however, these processors stop performing computations, and their outgoing messages may be dropped.

We consider algorithms that tolerate up to $t \leq \lceil n/2 \rceil - 1$ processor failures. That is, when more than half of the processors are non-faulty, they all return an answer from the protocol with probability one. A standard assumption in this setting is that all non-faulty processors always take part in the computation by replying to the messages, irrespective of whether they participate in a certain algorithm or even after they return a value—otherwise, the $t \leq \lceil n/2 \rceil - 1$ condition may be violated.

The Communicate Primitive. Our algorithms use a procedure called *communicate*, defined in [ABND95] as a building block for asynchronous communication. The call *communicate*(m) sends the message m to all n processors and waits for at least $\lfloor n/2 \rfloor + 1$ acknowledgments before proceeding with the protocol. The *communicate* procedure can be viewed as a best-effort broadcast mechanism; its key property is that any two *communicate* calls intersect in at least one recipient. In the following, a processor i will communicate messages of the form (*propagate*, v_i) or (*collect*, v). For the first message type, each recipient j updates its view of the variable v and acknowledges by sending back an *ACK* message. In the second case, the acknowledgement is a pair (*ACK*, v_j) containing j 's view of the variable for the receiving process. In both cases, processor i waits for $> n/2$ *ACK* replies before proceeding with its protocol. In the case of *collect*, the *communicate* call returns an array of at least $\lfloor n/2 \rfloor + 1$ views that were received.

Adversary. We consider strong adversarial setting where the scheduling of processor steps, message deliveries and processor failures are controlled by an adaptive adversary. At any point, this adversary can examine the state of the system, including the outcomes of random coin flips performed by the processors, and can adjust the scheduling accordingly.

Complexity Measures. We consider two worst-case complexity measures against the adaptive adversary: the maximum expected number of messages sent by all processors and the maximum expected number of *communicate* calls by any processor, that the adversary can achieve.

Problem Statements. In the *leader election* (*test-and-set*) problem, each processor may return either *WIN* or *LOSE*. A correct algorithm ensures that every (correct) processor returns (*termination*), and that only one processor may return *WIN* (*unique winner*). Also, no processor may lose before the eventual winner starts its execution. More precisely, the goal is to ensure that the processor operations are *linearizable*, i.e., can be ordered such (1) the first operation in the order is *WIN* and every other return value is *LOSE*, and (2) the order of non-overlapping operations is respected. *Strong (tight) renaming* requires each of the n participants to eventually return a name between 1 and n (*termination*), and each name must be *unique*.

3 The Leader Election Algorithm

Our leader election algorithm guarantees that if k processors participate, the maximum expected number of *communicate* calls by any processor that the *strong adaptive* adversary can achieve is $O(\log^* k)$, and the maximum expected total number of messages is $O(nk)$. We start by illustrating the main algorithmic idea.

3.1 The PoisonPill Technique

Consider the protocol specified in Figure 1 from the point of view of a participating processor. The procedure receives the id of the processor as an input, and returns a *SURVIVE/DIE* indication. All n processors react to received messages by replying with acknowledgments according to the *communicate* procedure. In the following, we call a *quorum* any set of more than $n/2$ processors.

Each participating processor announces that it is about to flip a random coin by moving to state *Commit* (lines 2-3), then obtain either low or high priority based on the outcome of a biased coin flip. The processor

Input: Unique identifier i of the participating processor
Output: *SURVIVE* or *DIE*
Local variables:
 $Status[n] = \{\perp\};$
 $Views[n][n];$
int $coin;$

```

1 procedure PoisonPill( $i$ )
2    $Status[i] \leftarrow Commit$                                 /* commit to coin flip */
3    $communicate(propagate, Status[i])$                         /* propagate status */
4    $coin \leftarrow random(1 \text{ with probability } 1/\sqrt{n}, 0 \text{ otherwise})$  /* flip coin */
5   if  $coin = 0$  then  $Status[i] \leftarrow Low-Pri$ 
6   else  $Status[i] \leftarrow High-Pri$ 
7    $communicate(propagate, Status[i])$                         /* propagate updated status */
8    $Views \leftarrow communicate(collect, Status)$            /* collect status from  $> n/2$  */
9   if  $Status[i] = Low-Pri$  then
10    if  $\exists proc. j : (\exists k : Views[k][j] \in \{Commit, High-Pri\} \text{ and } \forall k' : Views[k'][j] \neq Low-Pri)$  then
11      return DIE /*  $i$  has low priority, sees processor  $j$  with either high
12        priority or committed and not low priority, and dies */
13  return SURVIVE

```

Figure 1: PoisonPill Technique

then propagates its priority information to a quorum (line 7). Next, it collects the status of other processors from a quorum using the $communicate(collect, Status)$ call on line 8 that requests views of the array $Status$ from each processor j , returning the set of replies received, of size at least $n/2$.

The crux of the round procedure is the *DIE* condition on line 11. A processor p returns *DIE* at this line if *both* of the following occur: (1) the processor p has low priority, *and* (2) it observes another processor q that does not have low priority in any of the views, but q has either high priority or is committed to flipping a coin (state *Commit*) in some view. Otherwise, processor p survives. The first key observation is that

Claim 3.1. *If all processors participating in PoisonPill return, at least one processor survives.*

Proof. Assume the contrary. Since processors with high priority always survive, all participating processors must have a low priority. All participants propagate their low priority information to a quorum by calling the $communicate$ procedure on line 7. Let i be the last processor that completes this $communicate$ call. At this point, the status information of all participants is already propagated to a quorum. More precisely, for every participating processor j , more than half of the processors have a view $Status[j] = Low-Pri$.

Therefore, when processor i proceeds to line 8 and collects the $Status$ arrays from more than half of the processors, then, since any two quorums intersect, for every participating processor j , there will be a view of some processor k' showing j 's low priority. All non-participating processors will have priority \perp in all views. But given the structure of the protocol, processor i will *not* return on line 11 and will survive. This contradiction completes the proof. \square

On the other hand, we can bound the maximum expected number of processors that survive:

Claim 3.2. *The maximum expected number of processors that return *SURVIVE* is $O(\sqrt{n})$.*

Proof. Consider the random coin flips on line 4 and let us highlight the first time when some processor i flips value 1. We will argue that no other processor j that subsequently (or simultaneously) flips value 0 can survive. Consider such a state. When processor j flips 0, processor i has already propagated its *Commit* status to a quorum of processors. Furthermore, processor i has a high priority, thus no processor can ever view it as having a low priority. Hence, when processor j collects views from a quorum, because every two

quorums have an intersection, some processor k will definitely report the status of processor i as *Commit* or *High-Pri* and no processor will report *Low-Pri*. Thus, processor j will have to return *DIE* on line 11.

The above argument implies that processors survive either if they flip 1 and get a high priority, or if they flip 0 strictly before any other processor flips 1. Each of the at most n processors independently flips a biased coin on line 4 and hence, the number of processors that flip 1 is at most the number of 1's in n Bernoulli trials with success probability $1/\sqrt{n}$, in expectation \sqrt{n} . Processors that flip 0 at the same time as the first 1 do not survive, and it also takes \sqrt{n} trials in expectation before the first 1 is flipped, giving an upper bound \sqrt{n} on the maximum expected number of processors that can flip 0 and survive. \square

It is possible to apply this technique recursively with some extra care and construct an algorithm with an expected $O(\log \log n)$ round complexity. But we do not want to stop here.

3.2 Heterogeneous PoisonPill

Building a more efficient algorithm based on the PoisonPill technique requires reducing the number of survivors beyond $\Omega(\sqrt{n})$ without violating the invariant that not all participants may die. We control the coin flip bias, but setting the probability of flipping 1 to $1/\sqrt{n}$ is provably optimal. Let the adversary schedule processors to execute PoisonPill sequentially. With a larger probability of flipping 1, more than \sqrt{n} processors are expected to get a high priority and survive. With a smaller probability, at least the first \sqrt{n} processors are expected to all have low priority and survive. There are always $\Omega(\sqrt{n})$ survivors.

To overcome the above lower bound, after committing, we make each processor record the list ℓ of all processors including itself, that have a non- \perp status in some view collected from the quorum. Then we use the size of list ℓ of a processor to determine its probability bias. Each processor also augments priority with its ℓ and propagates that as a status. This way, every time a high or low priority of a processor p is observed, ℓ of processor p is also known. Finally, the survival criterion is modified: each processor first computes set L as the union of all processors whose non- \perp statuses it ever observed itself, and of the ℓ lists it has observed in priority informations in these statuses. If there is a processor in L for which no reported view has low priority, the current processor drops.

The algorithm is described in Figure 2. The particular choice of coin flip bias is influenced by factors that should become clear from the analysis. Despite modifications, the same argument as in Claim 3.1 still guarantees at least one survivor. Let us now prove that the views of the processors have the following interesting *closure property*, which will be critical to bounding the number of survivors with low priority.

Claim 3.3. *Consider any set S of processors that each flip 0 and survive. Let U be the union of all L lists of processors in S . Then, for $p \in U$ and every processor q in the ℓ list of p , q is also in U .*

Proof. In order for processors in S to survive, they should have observed a low priority for each of the processors in their L lists. Thus, every processor $p \in U$ must flip 0, as otherwise it would not have a low priority. However, the low priority of p observed by a survivor was augmented by the ℓ list of p . According to the algorithm, the survivor includes in its own L all processors q from this ℓ list of p , implying $q \in U$. \square

Next, let us prove a few other useful claims:

Claim 3.4. *If processor q completed executing line 15 no later than processor p completed executing line 15, then q will be included in the ℓ list of p .*

Proof. When p collects statuses on line 16 from a quorum, q is already done propagating its *Commit* on line 15. As every two quorum has an intersection, p will observe a non- \perp status of q on line 17. \square

Claim 3.5. *The probability of at least z processors flipping 0 and surviving is $O(1/z)$.*

```

13 procedure HeterogeneousPoisonPill( $i$ )
14    $Status[i] \leftarrow \{.stat = Commit, .list = \{\}\}$            /* commit to coin flip */
15    $communicate(propagate, Status[i])$                          /* propagate status */
16    $Views \leftarrow communicate(collect, Status)$              /* collect status from  $> n/2$  */
17    $\ell \leftarrow \{j \mid \exists k : Views[k][j] \neq \perp\}$        /* record participants */
18   if  $|\ell| = 1$  then  $prob \leftarrow 1$                        /* set bias */
19   else  $prob \leftarrow \frac{\log |\ell|}{|\ell|}$                        /* set bias */
20    $coin \leftarrow \text{random}(1 \text{ with probability } prob, 0 \text{ otherwise})$  /* flip coin */
21   if  $coin = 0$  then  $Status[i] \leftarrow \{.stat = Low-Pri, .list = \ell\}$  /* record priority and list */
22   else  $Status[i] \leftarrow \{.stat = High-Pri, .list = \ell\}$  /* record priority and list */
23    $communicate(propagate, Status[i])$                          /* propagate priority and list */
24    $Views \leftarrow communicate(collect, Status)$              /* collect status from  $> n/2$  */
25   if  $Status[i].stat = Low-Pri$  then
26      $L \leftarrow \cup_{k,j: Views[k][j] \neq \perp} Views[k][j].list$  /* union all observed lists */
27      $L \leftarrow L \cup \{j \mid \exists k : Views[k][j] \neq \perp\}$  /* record new participants */
28     if  $\exists proc. j \in L : \forall k : Views[k][j].stat \neq Low-Pri$  then
29       return DIE /*  $i$  has low priority, learns about processor  $j$ 
30         participating whose low priority is not reported, and dies */
31   return SURVIVE

```

Figure 2: Heterogeneous PoisonPill

Proof. Let S be the set of the z processors that flip 0 and survive and let us define U as in Claim 3.3. For any processor $p \in U$ and any processor q that completes executing line 15 no later than p , by Claim 3.4 processor q has to be contained in the ℓ list of p , which by the closure property (Claim 3.3) implies $q \in U$. Thus, if we consider the ordering of processors according to the time they complete executing line 15, all processors not in U must be ordered strictly after all processors in U .

Therefore, during the execution, first $|U|$ processors that complete line 15 must all flip 0. The adversary may influence the composition of U , but by the closure property, each ℓ list of processors in U contains only processors in U , meaning $|\ell| \leq |U|$. So the probability for each processor to flip 0 is at most $(1 - \frac{\log |U|}{|U|})$ and for all processors in U to flip 0's is at most $(1 - \frac{\log |U|}{|U|})^{|U|} = O(1/|U|)$. This is $O(1/z)$ since all z survivors from S are included in their own lists and hence also in U . \square

We have never relied on knowing n . If $k \leq n$ processors participate in the heterogeneous PoisonPill, we get

Lemma 3.6. *The maximum expected number of processors that flip 0 and survive is $O(\log k) + O(1)$.*

Proof. Let Z denote the number of processors that flip 0 and survive. Then, $\mathbb{E}[Z] = \sum_{z=0}^k \Pr[Z \geq z]$ and $\Pr[Z \geq z] = O(1/z)$ by Claim 3.5. \square

Lemma 3.7. *The maximum expected number of processors that flip 1 is $O(\log^2 k) + O(1)$.*

Proof. Consider the ordering of processors according to the time they complete executing line 15, breaking ties arbitrarily. Due to Claim 3.4, the processor that is ordered first always has $|\ell| \geq 1$, the second processor always computes $|\ell| \geq 2$, and so on. The probability of flipping 1 decreases as $|\ell|$ increases, and the best expectation achievable by adversary is $1 + \sum_{l=2}^k \frac{\log l}{l} = O(\log^2 k) + O(1)$ as desired. \square

3.3 Final construction

The idea of implementing leader election is to have rounds of heterogeneous PoisonPill, where all processors participate in the first round and only the survivors of round r participate in round $r + 1$. Each processor p , before participating in round r_p , first propagates r_p as its current round number to a quorum, then collects

information about the rounds of other processors from a quorum. Let R be the maximum round number of a processor in all views that p collected. To determine the winner, we use the idea from [SSW91]: if $R > r_p$, then p loses and if $R < r_p - 1$ then p wins. We also use a standard doorway mechanism [AGTV92] to ensure linearizability. The pseudocode of the final construction is given in Appendix A.1, along with the complete proof of the following statement:

Theorem A.5. *Our leader election algorithm is linearizable. If there are at most $\lceil n/2 \rceil - 1$ processor faults, all non-faulty processors terminate with probability 1. For k participants, the maximum expected number of communicate calls by any processor that the adaptive adversary can achieve is $O(\log^* k)$ and the maximum expected number of total messages is $O(kn)$.*

The performance guarantees follow from Lemma 3.6 and Lemma 3.7 with some careful analysis. In particular, later rounds in which maximum expected number of participants is constant require special treatment.

4 The Renaming Algorithm

```

Input: Unique identifier  $i$  from a large namespace
Output: int name  $\in [n]$ 
Local variables:
  bool Contended[n] = {false};
  bool Views[n][n];
  int coin, spot, outcome;
31 procedure getName( $i$ )
32   while true do
33     Views  $\leftarrow$  communicate(collect, Contended) /* collect contention information */
34     for  $j \leftarrow 1$  to  $n$  do
35       if  $\exists k : Views[k][j] = \text{true}$  then
36         Contended[j]  $\leftarrow$  true /* mark names that became contended */
37     communicate(propagate, { Contended[j] | Contended[j] = true }) /* propagate */
38     spot  $\leftarrow$  random( $j$  | Contended[j] = false) /* pick random uncontended name */
39     Contended[spot]  $\leftarrow$  true
40     outcome  $\leftarrow$  LeaderElectspot( $i$ ) /* contend for a new name */
41     communicate(propagate, Contended[spot]) /* propagate contention */
42     if outcome = WIN then
43       return spot /* win iff you are leader */

```

Figure 3: Pseudocode of the renaming algorithm for n processors.

The algorithm is described in Figure 3. There is a separate leader election protocol for each name; which a processor must win in order to claim the name. Each processor repeatedly chooses a new name and contends for it by participating in the corresponding leader election, until it eventually wins. Processors keep track of contended names and use this information to choose the next name to compete for: in particular, the next name is selected uniformly at random from the uncontended names. The algorithm is correct.

Lemma A.6. *No two processors return the same name from the getName call and if there are at most $\lceil n/2 \rceil - 1$ processor faults, all non-faulty processors terminate with probability 1.*

Omitted proofs can be found in Appendix A.2. Let us now introduce some notation. For an arbitrary execution, and for every name u , consider the first time when more than half of processors have `Contended[u] = true` in their view (or time ∞ , if this never happens). Let \prec denote the name ordering based on these times, and let $\{u_i\}$ be the sequence of names sorted according to increasing \prec . Among the names with time ∞ , sort later the ones that are never contended by the processors. Resolve all the remaining ties according to the order of the names. This ordering has the following useful temporal property.

Lemma A.7. *In any execution, if a processor views $\text{Contended}[i] = \text{true}$ in some while loop iteration, and in some subsequent iteration on line 38 the same processor views $\text{Contended}[j] = \text{false}$, $i \prec j$ has to hold.*

Let X_i be a random variable, denoting the number of processors that ever contend in a leader election for the name u_i . The following holds.

Lemma A.8. *The message complexity of our renaming algorithm is $O(n \cdot \mathbb{E}[\sum_{i=1}^n X_i])$.*

We partition names $\{u_i\}$ into $\log n$ groups, where the first group G_1 contains the first $n/2$ names, second group G_2 contains the next $n/4$ names, etc. We use notation $G_{j' \geq j}$, $G_{j'' > j}$ and $G_{j''' < j}$ to denote the union of all groups $G_{j'}$ where $j' \geq j$, all groups $G_{j''}$ where $j'' > j$, and all groups $G_{j'''}$ where $j''' < j$, respectively. We can now split any execution into at most $\log n$ phases. The first phase starts when the execution starts and ends as soon as for each $u_i \in G_1$ more than half of the processors view $\text{Contended}[u_i] = \text{true}$ (the way $\{u\}$ is sorted, this is the same as when the contention information about $u_{n/2}$ is propagated to a quorum). At this time, the second phase starts and ends when for each $u_i \in G_2$ more than half of the processors view $\text{Contended}[u_i] = \text{true}$. When the second phase ends, the third phase starts, and so on.

Consider any loop iteration of some processor p in some execution. We say that an iteration *starts* at a time instant when p executes line 32 and reaches line 33. Let V_p be p 's view of the *Contended* array right before picking a spot on line 38 in the given iteration. We say that an iteration is *clean*(j), if the iteration starts during phase j and no name from later groups $G_{j'' > j}$ is contended in V_p . We say that an iteration is *dirty*(j), if the iteration starts during phase j and some name from a later group $G_{j'' > j}$ is contended in V_p .

Observe that any iteration that starts in phase j can be uniquely classified as *clean*(j) or *dirty*(j) and in these iterations, processors view all names $u_i \in G_{j''' < j}$ from previous groups as contended.

Lemma A.9. *In any execution, at most $\frac{n}{2^j-1}$ processors ever contend for names from groups $G_{j' \geq j}$.*

Lemma 4.1. *For any fixed j , the total number of *clean*(j) iterations is larger than or equal to $\alpha n + \frac{n}{2^j-1}$ with probability at most $e^{-\frac{\alpha n}{16}}$ for all $\alpha \geq \frac{1}{2^j-5}$.*

Proof. Fix some j . Consider a time t when the first *dirty*(j) iteration is completed. At time t , there exists an i such that $u_i \in G_{j'' > j}$ and a quorum of processors view $\text{Contended}[i] = \text{true}$, so all iterations that start later will set $\text{Contended}[i] \leftarrow \text{true}$ on line 36. Therefore, any iteration that starts after t must observe $\text{Contended}[i] = \text{true}$ on line 38 and by definition cannot be *clean*(j). By Lemma A.9, at most $\frac{n}{2^j-1}$ processors can have active *clean*(j) iterations at time t . The total number of *clean*(j) iterations is thus upper bounded by $\frac{n}{2^j-1}$ plus the number of *clean*(j) iterations completed before time t , which we denote as *safe* iterations.² Safe iterations all finish before any iteration where a processor contends for a name in $G_{j'' > j}$ is completed. By Lemma A.9, at most $\frac{n}{2^j}$ different processors can ever contend for names in $G_{j'' > j}$, therefore, αn safe iterations can occur only if in at most $\frac{n}{2^j}$ of them processors choose to contend in $G_{j'' > j}$. Otherwise, some processor would have to complete an iteration where it contended for a name in $G_{j'' > j}$.

In every *clean*(j) iteration, on line 38, any processor p contends for a name in $G_{j' \geq j}$ uniformly at random among non-contended spots in its view V_p . With probability at least $\frac{1}{2}$, p contends for a name from $G_{j'' > j}$, because by definition of *clean*(j), all spots in $G_{j'' > j}$ are non-contended in V_p .

Let us describe the process by considering a random variable $Z \sim \text{B}(\alpha n, \frac{1}{2})$ for $\alpha \geq \frac{1}{2^j-5}$, where each success event corresponds to an iteration contending in $G_{j'' > j}$. By the Chernoff Bound, the probability of αn iterations with at most $\frac{n}{2^j}$ processors contending in $G_{j'' > j}$ is:

$$\Pr \left[Z \leq \frac{n}{2^j} \right] = \Pr \left[Z \leq \frac{\alpha n}{2} \left(1 - \frac{2^{j-1}\alpha - 1}{2^{j-1}\alpha} \right) \right] \leq \exp \left(-\frac{\alpha n (2^{j-1}\alpha - 1)^2}{2(2^{j-1}\alpha)^2} \right) \leq e^{-\frac{\alpha n}{8}}$$

So far, we have assumed that the set of names belonging to the later groups $G_{j'' > j}$ was fixed, but the adversary controls the execution. Luckily, what happened before phase j (i.e. the actual names that were acquired from $G_{j''' < j}$) is irrelevant, because all the names from the earlier phases are viewed as contended

² If no *dirty*(j) iteration ever completes, then we call all *clean*(j) iterations safe.

by all iterations that start in phases $j' \geq j$. Unfortunately, however, the adversary also influences what names belong to group j and to groups $G_{j'' > j}$. There are $\binom{2^{1-j}n}{2^{-j}n}$ different possible choices for names in G_j , and by a union bound, the probability that αn iterations can occur even for one of them is at most:

$$e^{-\frac{\alpha n}{8}} \cdot \binom{2^{1-j}n}{2^{-j}n} \leq e^{-\frac{\alpha n}{8}} \cdot (2e)^{2^{-j}n} \leq e^{-n(2^{-3}\alpha - 2^{1-j})} \leq e^{-\frac{\alpha n}{16}}$$

proving the claim. \square

Plugging $\alpha = \beta - \frac{1}{2^{j-1}} \geq \frac{1}{2^{j-5}}$ in the above lemma, we obtain the following

Corollary 4.2. *The total number of $\text{clean}(j)$ iterations is larger than or equal to βn with probability at most $e^{-\frac{\beta n}{32}}$ for all $\beta \geq \frac{1}{2^{j-6}}$.*

Let $X_i(\text{clean})$ be the number of processors that ever contend for $u_i \in G_j$ in some $\text{clean}(j)$ iteration and define $X_i(\text{dirty})$ analogously: as the number of processors that ever contend for $u_i \in G_j$ in some $\text{dirty}(j)$ iteration. Relying on Corollary 4.2, we get the following result:

Lemma A.10. $\mathbb{E}[\sum_{i=1}^n X_i(\text{clean})] = O(n)$.

Each iteration where a processor contends for a name $u_i \in G_j$ is by definition either as $\text{clean}(j)$, $\text{dirty}(j)$ or starts in a phase $j''' < j$. Let us call these $\text{cross}(j)$ and denote by $X_i(\text{cross})$ the number of such iterations. We show that in any execution, for each j , any processor participates in at most one $\text{dirty}(j)$ and at most one $\text{cross}(j)$ iteration. This allows us to prove with some work that $\mathbb{E}[\sum_{i=1}^n X_i(\text{dirty})] = O(n)$ and $\mathbb{E}[\sum_{i=1}^n X_i(\text{cross})] = O(n)$ (Lemma A.12). The message complexity upper bound then follows by piecing together the previous claims.

Theorem 4.3. *The expected message complexity of our renaming algorithm is $O(n^2)$.*

Proof. We know $X_i = X_i(\text{clean}) + X_i(\text{dirty}) + X_i(\text{cross})$ and by Lemma A.10, Lemma A.12 we get $\mathbb{E}[\sum_i X_i] = O(n)$. Combining with Lemma A.8 gives the desired result. \square

The round complexity upper bound exploits a trade-off between the probability that a processor collides in an iteration (and must continue) and the ratio of available slots which must be assigned during that iteration.

Theorem A.13. *In the renaming algorithm, the maximum expected number of communicate calls by any processor that the adaptive adversary can achieve is $O(\log^2 n)$.*

5 Message Complexity Lower Bounds

We can prove that the algorithms we presented for leader election and renaming in the previous two sections are asymptotically message-optimal. Due to space constraints, the proof of this result is deferred to the Appendix. (In fact, we show a stronger claim, proving the same message complexity lower bound for arbitrary objects with strongly non-commutative operations [AGH⁺11].)

Corollary B.3. *Any implementation of leader election or renaming by $k \leq n$ processors guaranteeing termination with probability at least $\alpha > 0$ in an asynchronous message-passing system where $t < n/2$ processors may fail by crashing must have worst-case expected message complexity $\Omega(\alpha k n)$.*

6 Discussion and Future Work

We have given the first sub-logarithmic leader election algorithm against a strong adversary, and asymptotically tight bounds for the message complexity of renaming and leader election. Our results also limit the power of topological lower bound techniques, e.g. [HS99], when applied to randomized leader election, since these techniques allow processors to communicate using unit-cost broadcasts or snapshots. Our algorithm shows that no bound stronger than $\Omega(\log^* n)$ time is possible using such techniques, unless the cost of information dissemination is explicitly taken into account.

Determining the tight time complexity bounds for leader election remains an intriguing open question. Another interesting research direction would be to apply the tools we developed to obtain time- and message-efficient implementations of other fundamental distributed tasks, such as task allocation or mutual exclusion, and to explore solutions optimizing bit complexity.

References

- [AA11] Dan Alistarh and James Aspnes. Sub-logarithmic test-and-set against a weak adversary. In *Distributed Computing: 25th International Symposium, DISC 2011*, volume 6950 of *Lecture Notes in Computer Science*, pages 97–109. Springer-Verlag, September 2011.
- [AACH⁺13] Dan Alistarh, James Aspnes, Keren Censor-Hillel, Seth Gilbert, and Rachid Guerraoui. Tight bounds for asynchronous renaming. Accepted to JACM, September 2013.
- [AAG⁺10] Dan Alistarh, Hagit Attiya, Seth Gilbert, Andrei Giurgiu, and Rachid Guerraoui. Fast Randomized Test-and-Set and Renaming. In *Proceedings of DISC 2010*, Lecture Notes in Computer Science. Springer-Verlag New York, Ms Ingrid Cunningham, 175 Fifth Ave, New York, Ny 10010 Usa, 2010.
- [AAKS14] Dan Alistarh, James Aspnes, Valerie King, and Jared Saia. Message passing meets shared memory: A new approach for communication efficient randomized consensus. In *Proceedings of the 28th International Symposium on Distributed Computing, DISC 2014*, February 2014.
- [ABGG12] Dan Alistarh, Michael A. Bender, Seth Gilbert, and Rachid Guerraoui. How to allocate tasks asynchronously. In *53rd Annual IEEE Symposium on Foundations of Computer Science, FOCS 2012, New Brunswick, NJ, USA, October 20-23, 2012*, pages 331–340, 2012.
- [ABND⁺90] Hagit Attiya, Amotz Bar-Noy, Danny Dolev, David Peleg, and Rüdiger Reischuk. Renaming in an asynchronous environment. *J. ACM*, 37(3):524–548, July 1990.
- [ABND95] Hagit Attiya, Amotz Bar-Noy, and Danny Dolev. Sharing memory robustly in message-passing systems. *Journal of the ACM (JACM)*, 42(1):124–142, 1995.
- [AC08] Hagit Attiya and Keren Censor. Tight bounds for asynchronous randomized consensus. *J. ACM*, 55(5):20:1–20:26, November 2008.
- [AGH⁺11] Hagit Attiya, Rachid Guerraoui, Danny Hendler, Petr Kuznetsov, Maged M Michael, and Martin Vechev. Laws of order: expensive synchronization in concurrent algorithms cannot be eliminated. In *ACM SIGPLAN Notices*, volume 46, pages 487–498. ACM, 2011.
- [AGTV92] Yehuda Afek, Eli Gafni, John Tromp, and Paul M. B. Vitányi. Wait-free test-and-set (extended abstract). In *Proceedings of the 6th International Workshop on Distributed Algorithms, WDAG '92*, pages 85–94, London, UK, UK, 1992. Springer-Verlag.
- [AKS83] M. Ajtai, J. Komlós, and E. Szemerédi. An $O(n \log n)$ sorting network. In *Proceedings of the Fifteenth Annual ACM Symposium on Theory of Computing, STOC '83*, pages 1–9, New York, NY, USA, 1983. ACM.
- [BKRS96] Jonathan F. Buss, Paris C. Kanellakis, Prabhakar L. Ragde, and Alex Allister Shvartsman. Parallel algorithms with processor failures and delays. *J. Algorithms*, 20:45–86, January 1996.
- [BO83] Michael Ben-Or. Another advantage of free choice (extended abstract): Completely asynchronous agreement protocols. In *Proceedings of the Second Annual ACM Symposium on Principles of Distributed Computing, PODC '83*, pages 27–30, New York, NY, USA, 1983. ACM.
- [Dij65] E. W. Dijkstra. Solution of a problem in concurrent programming control. *Commun. ACM*, 8(9):569–, September 1965.
- [FLP85] Michael J. Fischer, Nancy A. Lynch, and Michael S. Paterson. Impossibility of distributed consensus with one faulty process. *J. ACM*, 32(2):374–382, April 1985.
- [GW12a] George Giakkoupis and Philipp Woelfel. On the time and space complexity of randomized test-and-set. In *Proceedings of the 2012 ACM Symposium on Principles of Distributed Computing, PODC '12*, pages 19–28, New York, NY, USA, 2012. ACM.
- [GW12b] George Giakkoupis and Philipp Woelfel. A tight rnr lower bound for randomized mutual exclusion. In *Proceedings of the Forty-fourth Annual ACM Symposium on Theory of Computing, STOC '12*, pages 983–1002, New York, NY, USA, 2012. ACM.
- [Her91] Maurice Herlihy. Wait-free synchronization. *ACM Trans. Program. Lang. Syst.*, 13(1):124–149, January 1991.
- [HS99] Maurice Herlihy and Nir Shavit. The topological structure of asynchronous computability. *Journal of The ACM*, 46:858–923, 1999.

- [HW09] Danny Hendler and Philipp Woelfel. Randomized mutual exclusion in $o(\log n / \log \log n)$ rmrs. In *Proceedings of the 28th ACM Symposium on Principles of Distributed Computing*, PODC '09, pages 26–35, New York, NY, USA, 2009. ACM.
- [HW10] Danny Hendler and Philipp Woelfel. Adaptive randomized mutual exclusion in sub-logarithmic expected time. In *Proceedings of the 29th ACM SIGACT-SIGOPS Symposium on Principles of Distributed Computing*, PODC '10, pages 141–150, New York, NY, USA, 2010. ACM.
- [KS92] Paris C. Kanellakis and Alex A. Shvartsman. Efficient parallel algorithms can be made robust. *Distrib. Comput.*, 5(4):201–217, April 1992.
- [LSP82] Leslie Lamport, Robert Shostak, and Marshall Pease. The byzantine generals problem. *ACM Trans. Program. Lang. Syst.*, 4(3):382–401, July 1982.
- [Lyn97] Nancy Lynch. *Distributed Algorithms*. Morgan Kaufmann, 1997.
- [PSL80] M. Pease, R. Shostak, and L. Lamport. Reaching agreement in the presence of faults. *J. ACM*, 27(2):228–234, April 1980.
- [SSW91] Michael Saks, Nir Shavit, and Heather Woll. Optimal time randomized consensusmaking resilient algorithms fast in practice. In *Proceedings of the second annual ACM-SIAM symposium on Discrete algorithms*, pages 351–362. Society for Industrial and Applied Mathematics, 1991.

A Deferred Proofs

A.1 Leader Election Construction and Analysis

Figure 4 contains the pseudocode of PreRound procedure that processors execute before participating in round r . Every processor starts in the same initial non-negative round. The PreRound procedure takes round number r and the id of the processor as an input and outputs either *PROCEED*, *WIN* or *LOSE*. Each processor p first propagates r to a quorum, then collects information about the rounds of other processors also from a quorum. Let R be the maximum round number of a processor in all views that p collected. Using idea from [SSW91], if $R > r$, then p loses, if $R < r - 1$ then p wins and otherwise p returns *PROCEED*.

Input: Unique identifier i of the participating processor, round number r
Output: *PROCEED*, *WIN*, or *LOSE*
Local variables:

```

int Round[ $n$ ] = {0};
int Views[ $n$ ][ $n$ ];
int R;
44 procedure PreRound( $i, r$ )
45   Round[ $i$ ]  $\leftarrow r$                                /* record own round */
46   communicate(propagate, Round[ $i$ ])                 /* propagate own round */
47   Views  $\leftarrow$  communicate(collect, Round)         /* collect round from  $> n/2$  */
48    $R \leftarrow \max_{k, j | j \neq i} (\text{Views}[k][j])$  /* maximum round of other processors observed */
49   if  $r < R$  then
50     return LOSE
51   if  $R < r - 1$  then
52     return WIN
53   return PROCEED

```

Figure 4: PreRound procedure

To ensure linearizability we use a standard doorway technique, described in Figure 5. This doorway mechanism is implemented by the variable *door* stored by the processors. A value false corresponds to the door being open and a value true corresponds to the door being closed. Each participating processor p starts by collecting the views of *door* from more than half of the processors on line 56. If a closed door is reported, p is too late and automatically returns *LOSE*. The door is closed by processors on line 59, and

```

Output: PROCEED or LOSE
Local variables:
    bool door = false                                /* door is initially open */
54  bool Doors[n];
55  procedure Doorway( $\langle \rangle$ )
56    Doors  $\leftarrow$  communicate(collect, door)          /* collect door from  $> n/2$  */
57    if  $\exists j : \text{Doors}[j] = \text{true}$  then
58      return LOSE                                       /* lose if the door is closed */
59    door  $\leftarrow$  true                                   /* close the door */
60    communicate(propagate, door)                       /* propagates door = true to  $> n/2$  */
61    return PROCEED

```

Figure 5: Doorway procedure

this information is then propagated to a quorum. The goal of the doorway is to ensure that no processor can lose *before* the winner has started its execution.

Finally we put the pieces together. Our complete leader election algorithm is described in Figure 6. It involves going through the doorway procedure in the beginning, and then rounds of PreRound procedure possibly followed by participation in a HeterogeneousPoisonPill protocol for round r . Note that HeterogeneousPoisonPill protocols for different rounds are completely disjoint from each other.

```

Input: Unique identifier  $i$  of the participating processor
Output: WIN or LOSE
Local variables:
    int  $r = 1$ ;
    outcome;
62  procedure LeaderElect( $i$ )
63    if Doorway( $\langle \rangle$ ) = LOSE then
64      return LOSE                                       /* lose if door was closed */
65    repeat
66      outcome  $\leftarrow$  PreRound( $i, r$ )                    /* preround routine */
67      if outcome  $\in$  {WIN, LOSE} then
68        return outcome                                   /* return if rounds permit */
69      if HeterogeneousPoisonPill $_r$ ( $i$ ) = DIE then
70        return LOSE                                       /* lose if did not survive the round */
71       $r \leftarrow r + 1$ 
72    until never

```

Figure 6: Leader election algorithm

We now prove the properties of the algorithm.

Lemma A.1. *If all processors that call LeaderElect return, at least one processor returns WIN.*

Proof. Assume for contradiction that all processors that participate in PoisonPill return *LOSE*. Let us first prove that at least one processor always reaches the loop on line 65, or alternatively that not all processors can lose on line 64. This would mean that all processors return *LOSE* on line 58 of the Doorway procedure, but in that case the door would never be closed on line 59. Thus, all processor views would be $door = \text{false}$, and no processor would actually be able to return on line 58.

Since we showed that at least one processor reaches the loop, let us consider the largest round r in which some processors return, either in the pre-round routine of round r on line 68 or because of the poison pill

on line 70. By our assumption all these processors return *LOSE* in round r . But then, none of them may return on line 68, because this is only possible after returning *LOSE* on line 50, which only happens if a larger round than r is reported, contradicting our assumption that r is the largest round.

So, at least one processor participates in the *HeterogeneousPoisonPill_r* protocol. However, by exactly the same argument as in Claim 3.1, *HeterogeneousPoisonPill_r* is guaranteed to have at least one survivor which would then participate in round $r + 1$, again contradicting that r is the largest round. \square

Lemma A.2. *At most one processor that executes LeaderElect can return WIN.*

Proof. A processor p can only return *WIN* from LeaderElect on line 68, which only happens after p returns *WIN* from PreRound call with some round r . This means p first propagated round r to a quorum on line 46, then collected views of *Round* array on line 47, and observed maximum round $R < r - 1$ of any processor in any of the views. This implies that when p finished propagating r to a quorum, no processor had finished propagating $r - 1$, i.e. executing line 46 in round $r - 1$. Otherwise, since every two quorums have an intersection, p would have observed round $r - 1$ and $R < r - 1$ would not hold. But for every other processor q , when q executes line 47 in round $r - 1$ and invokes the PreRound procedure, R will be at least r since p has already propagated to a quorum, so q will observe $r - 1 < r$ and return *LOSE* on line 50 and subsequently return *LOSE* from LeaderElect. \square

Lemma A.3. *Our leader election algorithm is linearizable.*

Proof. All processors that execute LeaderElect cannot return *LOSE* by Lemma A.1. Therefore, in every execution we can find LeaderElect invocation where processor either does not return, or returns *WIN*. On the other hand, by Lemma A.2, no more than one processor can return *WIN*. If no processor returns *WIN*, let us linearize the processor that invoked LeaderElect the earliest as the leader. This way, we always have a unique processor to be linearized as the winner. We linearize it at the beginning of its invocation interval, say point P , and claim that every remaining LeaderElect call can be linearized as returning *LOSE* after P .

Assume contrary, then the problematic LeaderElect invocation must return before P , and we know it has to return *LOSE*. By definition, this earlier call either closes the door or observes a closed door while executing the Doorway procedure. Therefore, the later call that we are linearizing as the winner has to observe a closed door on line 56 and cannot avoid returning *LOSE* on line 58. Hence, this invocation can never return *WIN*, and since we are linearizing it as winner, it should be the case that it does not return and no other processor returns *WIN*. We picked this invocation to have the earliest starting point, so every other LeaderElect invocation that does not return must start after P . Let us now consider an extension of the current execution where the processors executing these invocations are continuously scheduled to take steps and all messages are delivered. According to the above argument, since all invocations start after P , these processors must observe a closed door on line 56 and return *LOSE* after only finitely many steps. We have hence constructed a valid execution where all processors that execute LeaderElect return *LOSE*. This contradiction with Lemma A.1 completes the proof. \square

We need one final claim before proving the main theorem.

Claim A.4. *The maximum expected number of participants decreases at least by some fixed constant fraction in every two rounds.*

Proof. This obviously holds for a single participant, because it will return *WIN* in the next round and the number of participants after that will be zero.

We know that for k participants in some round, by Lemma 3.6 and Lemma 3.7, the maximum expected number of participants in the next round is $O(\log^2 k + 1)$. This implies that for a large enough constant D , there is constant $c_1 < 1$ such that for $k > D$ the maximum expected number of participants in the next round, and thus in all rounds thereafter, is at most $c_1 k$. If $k \leq D$, then the first processor that finishes executing line 15 flips 1 with at least a constant probability. In this case, all processors that flip 0 will die,

and the expected number of the remaining processors that flip 0 is at least $\frac{k-1}{2} \leq \frac{k}{4}$ for $k \geq 2$. This is because the expected number of remaining processors that flip 1 is at most $\frac{k-1}{2}$, as each of them observes at least the first processor and itself, hence has no more than $1/2$ probability of flipping 1. Thus, if $k \leq D$, with a constant probability, a constant fraction of participants dies, meaning that there is a constant $c_2 < 1$ such that the maximum expected number of participants is at most $c_2 k$. Setting $c = \max(c_1, c_2) < 1$ we obtain that the maximum expected number of participants in every two rounds always decreases by at least a constant fraction to ck . \square

Theorem A.5. *Our leader election algorithm is linearizable. If there are at most $\lceil n/2 \rceil - 1$ processor faults, all non-faulty processors terminate with probability 1. For k participants, the maximum expected number of communicate calls by any processor that the adaptive adversary can achieve is $O(\log^* k)$ and the maximum expected number of total messages is $O(kn)$.*

Proof. We have shown linearizability in Lemma A.3.

All $k \geq 1$ processors participate in the first round. The maximum expected number of processors that participate in round 3 is clearly no more than the maximum expected number of survivors of the first round, which by Lemma 3.6 and Lemma 3.7 for $k > 1$ can be written as $C(\log^2 k + 2 \log k)$ for some constant C . If $k = 1$, then this lone processor will observe all other processors in round 0, leading to $R = 0$ and as current round is $r = 2$ it will return *WIN* in the second round. Hence, for $k = 1$, there will be zero participants in the third round. Thus, for any k , the maximum expected number of participants in round 3 is at most $f(k) = C(\log^2 k + 2 \log k)$.

Let us say the adversary can achieve a probability distribution for round 3 such that there are K_i participants with probability p_i . We have shown above that

$$\sum_i p_i K_i \leq f(k) \tag{A.1}$$

Now, using the same argument as above, we can bound the maximum expected number of participants in round 5 to be at most $\sum p_i f(K_i)$. Function f is concave for non-negative arguments, and for arguments larger than a constant it is also monotonically increasing. This implies that either $\sum p_i K_i$, the expected number of participants in round 3, is constant, or

$$\sum p_i f(K_i) \leq f(\sum p_i K_i) \leq f(f(k)) \tag{A.2}$$

where the first part is Jensen's inequality and the second follows from (A.1) and the monotonicity property. Similarly, we get that unless the maximum expected number of participants in round 5 is less than a constant, the maximum number of participants in round 7 is at most $f(f(f(k)))$, and so on. Since $f(f(k)) \geq \log k$ for all k larger than some constant, if we denote by S_0 the number of participants in round $1 + 2 \log^* k$, maximum $\mathbb{E}[S_0]$ that the adversary can achieve must also be constant. These S_0 participants execute the same algorithm, with S_1 of them participating in the next round, etc.

Let R be the number of remaining rounds. Expectation of R can be written as

$$\mathbb{E}[R] = \sum_{i=1}^{\infty} \Pr[R \geq i] = \sum_{i=1}^{\infty} \Pr[S_i \geq 1] \leq \sum_{i=1}^{\infty} \mathbb{E}[S_i] \tag{A.3}$$

where the equality is by the definition of rounds and then we apply Markov's inequality to get to expectations. Finally, by Claim A.4 we get that $\mathbb{E}[R] = O(\mathbb{E}[S_0]) = O(1)$ and thus the maximum total number of rounds any processor participates in is $O(\log^* k)$, and processors perform only fixed, constantly many communicate calls per round.

To bound the maximum expected number of messages, let Q_r be the number of participants in round r , counting from the very first round. Since each processor sends $O(n)$ messages per round, the maximum expected number of messages is $\sum_{r=1}^{\infty} \mathbb{E}[O(nQ_r)] = n \cdot \mathbb{E}[O(Q_1)] = O(nk)$ using Claim A.4.

If there are at most $\lceil n/2 \rceil - 1$ processor faults, all communicate calls return, and processors must enter larger rounds. However, the probability that all processors terminate before reaching round r is $1 - \Pr[Q_r \geq 1] \geq 1 - \mathbb{E}[Q_r]$ which tends to 1 as r increases by Claim A.4. \square

A.2 Renaming Analysis

Lemma A.6. *No two processors return the same name from the getName call and if there are at most $\lceil n/2 \rceil - 1$ processor faults, all non-faulty processors terminate with probability 1.*

Proof. Assume that less than half of the processors are faulty. Processors executing the getName call the communicate procedure which always terminates under at most at most $\lceil n/2 \rceil - 1$ processor faults. All local computations steps are also always performed successfully by non-faulty processors.

Finally, processors invoke our leader election algorithm from Section 3 for at most n names, at most once for each name (the first time they set $Contended \leftarrow \text{true}$, which prohibits contending in the future). By Theorem A.5 all invocations of the leader election for a particular name by non-faulty processors terminate with probability 1, and using union bound for at most n names, the probability that all leader election calls by all non-faulty processors terminate tends to 1. Therefore, with probability 1, non-faulty processors keep making progress, i.e. they keep contending for new names, and as there are n names and n processors that do not contend for the same name twice, each non-faulty processor eventually wins a name and returns.

A processor that returns some name u from a getName call has to be the winner of our leader election protocol. However, according to Theorem A.5, LeaderElect_u cannot have more than one winner. \square

Lemma A.7. *In any execution, if a processor views $Contended[i] = \text{true}$ in some while loop iteration, and in some subsequent iteration on line 38 the same processor views $Contended[j] = \text{false}$, $i \prec j$ has to hold.*

Proof. Clearly, $j \neq i$ because contention information never disappears from a processor's view. In the earlier iteration, the processor propagates $Contended[i] = \text{true}$ to a quorum on line 37 or 41. During the later iteration, on line 33, the processor collects information and does not set $Contended[j]$ to true before reaching line 38. Thus, more than half of the processors view $Contended[j] = \text{false}$ at some intermediate time point. Therefore, a quorum of processors views $Contended[i] = \text{true}$ strictly earlier than $Contended[j] = \text{true}$, and by definition $i \prec j$. \square

Lemma A.8. *The message complexity of our renaming algorithm is $O(n \cdot \mathbb{E}[\sum_{i=1}^n X_i])$.*

Proof. Let L_j be the number of loop iterations executed by processor j . Then $\sum_i X_i = \sum_j L_j$, because every iteration involves one processor contending at a single name spot, and no processor contends for the same name twice. Each iteration involves two communicate calls with $O(n)$ total messages. The total number of messages sent in the leader election protocols is $O(\sum_i n \cdot X_i)$. The message complexity is thus the expectation of:

$$O\left(\sum_i n \cdot X_i\right) + \sum_j O(n) \cdot L_j = O\left(\sum_i n \cdot X_i\right)$$

as desired. \square

Lemma A.9. *In any execution, at most $\frac{n}{2^{j-1}}$ processors ever contend for names from groups $G_{j' \geq j}$.*

Proof. If no name from $G_{j' \geq j}$ is ever contended, then the statement is trivially true. If some name from $G_{j' \geq j}$ was contended, then by our ordering so were all names from the former groups. Otherwise, an uncontended name from an earlier group must be sorted later and cannot belong to an earlier group.

There are $n - \frac{n}{2^{j-1}}$ names in earlier groups. Since they all were contended, there are $n - \frac{n}{2^{j-1}}$ processors that can be linearized to win the corresponding leader election and the name. Consider one such processor p and the name u from some earlier group $G_{j'' < j}$, that p is bound to win. Processor p does not contend

for names after u , and it also never contends for a name from $G_{j' \geq j}$ before contending for u , because that contradicts Lemma A.7. Thus, none of the $n - \frac{n}{2^{j-1}}$ processors ever contend for a name from $G_{j' \geq j}$, out of n processors in total, completing the argument. \square

Lemma A.10. $\mathbb{E}[\sum_{i=1}^n X_i(\text{clean})] = O(n)$.

Proof. Let us equivalently prove that $\mathbb{E}[X_i(\text{clean})] = O(1)$ for any name u_i in some group G_j , where $X_i(\text{clean})$ is defined as the number of $\text{clean}(j)$ iterations, in which processors contend for a name $u_i \in G_j$.

By definition, in all $\text{clean}(j)$ iterations a processor observes all names in $G_{j'' > j}$ as uncontended on line 38. Therefore, each time, independent of other iterations, the probability of picking spot i and contending for the name u_i is at most $\frac{2^j}{n}$. Thus, if there are exactly βn of $\text{clean}(j)$ iterations, $X_i(\text{clean}) \leq B(\beta n, \frac{2^j}{n})$, thus

$$\mathbb{E}[X_i(\text{clean}) \mid u_i \in G_j, \beta n \text{ iterations}] \leq 2^j \beta \quad (\text{A.4})$$

for βn clean iterations that started in phase j . The probability that there are exactly βn of $\text{clean}(j)$ iterations is trivially upper-bounded by the probability that there are at least βn $\text{clean}(j)$ iterations, which by Corollary 4.2 is at most $e^{-\frac{\beta n}{32}}$ for $\beta \geq \frac{1}{2^{j-6}}$. Therefore:

$$\mathbb{E}[X_i(\text{clean}) \mid u_i \in G_j] \leq \frac{2^j}{2^{j-6}} + \sum_{l=\lceil \frac{n}{2^{j-6}} \rceil}^{\infty} e^{-\frac{l}{32}} \cdot \frac{2^j l}{n} \quad (\text{A.5})$$

which, after some calculation, is $O(1)$, completing the proof. \square

Claim A.11. *In any execution, for each j , any processor participates in at most one $\text{dirty}(j)$ and at most one $\text{cross}(j)$ iteration.*

Proof. The first time processor p participates in a $\text{dirty}(j)$ iteration, by definition, it views $\text{Contended}[i] = \text{true}$ for some $u_i \in G_{j'' > j}$. Therefore, p also propagates $\text{Contended}[i] = \text{true}$ on line 37 in the same iteration. When p starts a subsequent iteration, a quorum of processors know about $u_i \in G_{j'' > j}$ being contended. By the way names in u are sorted, at that point more than half of the processors must already know that each name in G_j is contended, meaning that phase j has ended. Therefore, no subsequent iteration of the processor can be of a $\text{dirty}(j)$ type.

On the other hand, when a processor completes a $\text{cross}(j)$ iteration, it has propagated contention information for $u_i \in G_j$ to a quorum, meaning that because of the way names in u are sorted, phase j must have been started already, and no operation that starts later can be $\text{cross}(j)$. \square

Lemma A.12. $\mathbb{E}[\sum_{i=1}^n X_i(\text{dirty})] = O(n)$ and $\mathbb{E}[\sum_{i=1}^n X_i(\text{cross})] = O(n)$.

Proof. Recall that $X_i(\text{dirty})$ is the number of processors that ever contend for $u_i \in G_j$ in a $\text{dirty}(j)$ iteration. Let us equivalently fix j and prove that $\mathbb{E}[\sum_{u_i \in G_j} X_i(\text{dirty})] = O(\frac{n}{2^{j-1}})$, which implies the desired statement by linearity of expectation and telescoping.

We sum up quantities $X_i(\text{dirty})$ for the names in j -th group, but the adversary controls precisely which names belong to group G_j . We will therefore consider all names $u_i \in G_{j' \geq j}$ and sum up quantities $X_{i,j}$: the number of processors that contend for a name u_i in a $\text{dirty}(j)$ iteration. All $\text{dirty}(j)$ iterations by definition start in phase j , and by Lemma A.9 there can be at most $\frac{n}{2^{j-1}}$ different processors executing them. Moreover, by Claim A.11 each of these processors can participate in at most one $\text{dirty}(j)$ iteration, implying $\sum_{u_i} X_{i,j} \leq \frac{n}{2^{j-1}}$. Thus $\mathbb{E}[\sum_{i=1}^n X_i(\text{dirty})] = \sum_{j=1}^{\log n} \sum_{u_i \in G_{j' \geq j}} X_{i,j} = O(n)$ as desired.

The proof for $X_i(\text{cross})$ is analogous because at most $\frac{n}{2^{j-1}}$ different processors contend for names $u_i \in G_{j' \geq j}$ by Lemma A.9, each participating in at most one $\text{cross}(j)$ iteration by Claim A.11. \square

Theorem A.13. *In the renaming algorithm, the maximum expected number of communicate calls by any processor that the adaptive adversary can achieve is $O(\log^2 n)$.*

Proof. In the following, we fix an arbitrary processor p , and upper bound the number of loop iterations it performs during the execution. Let M_i be the set of free slots that p sees when performing its random choice in the i th iteration of the loop, and let $m_i = |M_i|$. By construction, notice that there can be at most m_i processors that compete with for slots in M_i for the rest of the execution.

Assuming that p does not complete in iteration i , let $Y_i \subseteq M_i$ be the set of *new* slots that p finds out have become contended at the beginning of iteration $i + 1$, and let $y_i = |Y_i|$. We define an iteration as being *low-information* if $y_i/m_i < 1/\log m_i$. Notice that, in an iteration that is high-information, the processor might collide, but at least reduces its random range for choices by a $1/\log m_i$ factor.

Let us now focus on low-information iterations, and in particular let i be such an iteration. Notice that we can model the interaction between the algorithm and the adversary in iteration i as follows. Processor p first makes a random choice r from m_i slots it sees as available. By the principle of deferred decisions, we can assume that, at this point, the adversary schedules all other $m_i - 1$ processors to make their choices in this round, from slots in M_i , with the goal of causing a collision with p 's choice. (The adversary has no interest in showing slots outside M_i to processors.) Notice that, in fact, the adversary may choose to schedule certain processors multiple times in order to obtain collisions. However, by construction, each re-scheduled processor announces its choice in the iteration to a quorum, and this choice will become known to p in the next iteration. Therefore, re-scheduled processors should not announce more than $m_i/\log m_i$ distinct slots. Intuitively, the number of re-schedulings for the adversary can be upper bounded by the number of balls falling into the $m_i/\log m_i$ most loaded bins in an $m_i - 1$ balls into m_i bins scenario. A simple balls-and-bins argument yields that, in any case, the adversary cannot perform more than m_i re-schedules without having to announce $m_i/\log m_i$ new slots, with high probability in m_i .

Recall that the goal of the adversary is to cause a collision with p 's random choice r . We can reduce this to a balls-into-bins game in which the adversary throws $m_i - 1$ initial balls and an extra m_i balls (from the re-scheduling) into a set of $m_i(1 - 1/\log m_i)$ bins, with the goal of hitting a specific bin, corresponding to r . (The extra $(1 - 1/\log m_i)$ factor comes from the fact that certain processors (or balls) may already observe the slots removed in this iteration.) The probability that a fixed bin gets hit is at most

$$\left(1 - \frac{1}{m_i(1 - 1/\log m_i)}\right)^{2m_i} \leq (1/e)^3.$$

Therefore, processor p terminates in each low-information iteration with constant probability. Putting it all together, we obtain that, for $c \geq 4$ constant, after $c \log^2 n / \log \log n$ iterations, any processor p will terminate with high probability, either because $m_i = 1$ or because one of its probes was successful in a low-information phase.

In each loop iteration, a processor performs a fixed constant additional number of communicate calls on top of the communicate calls performed while executing the leader election algorithm for the name picked in that iteration. By Theorem A.5, the maximum expected number of communicate calls in each leader election is $O(\log^* n)$, and by linearity of expectation, total maximum number of communicate calls by any processor is at most $O(\frac{\log^2 n \log^* n}{\log \log n}) = O(\log^2 n)$. \square

B Message Complexity Lower Bounds

In this section, we prove that our algorithms are message-optimal by showing a lower bound of expected $\Omega(n^2)$ messages on any algorithm implementing leader election or renaming in an asynchronous message-passing system where $t < n/2$ processors may fail by crashing. In fact, we prove such a lower bound for any object with *strongly non-commutative* methods [AGH⁺11].

Definition B.1. Given an object O , a method M of this object is strongly non-commutative if there exists some state S of O for which an instance m_1 of M executed sequentially by processor p changes the result of an instance m_2 of M executed by processor $q \neq p$, and vice-versa, i.e. m_2 changes the result of m_1 from state S .

We now give a message complexity lower bound for objects with non-commutative operations.

Theorem B.2. Any implementation of an object O with a strongly non-commutative operation M by $k \leq n$ processors guaranteeing termination with at least constant probability $\alpha > 0$ in an asynchronous message-passing system where $t < n/2$ processors may fail by crashing must have worst-case expected message complexity $\Omega(\alpha kn)$.

Proof. Let A be an algorithm implementing a shared object O with a strongly non-commutative method M , in asynchronous message-passing with $t < n/2$, guaranteeing termination with probability α . We define an adversarial strategy for which we will argue that all the resulting terminating executions (regardless of their probability) must cause $\Omega(kn)$ messages to be sent. This clearly implies our claim. The strategy proceeds as follows.

Assume that each processor is executing an instance of M . The adversary picks a subset S of $k/4$ participants, and places them in a “bubble:” for each such processor q , the adversary suspends all its incoming and outgoing messages in a buffer, until there are at least $n/4$ such messages in the buffer. At this point, the processor is freed from the bubble, and is allowed to take steps synchronously, together with other processors. Processors outside S execute in lock-step, and their messages outside the bubble are delivered in a timely fashion.

Note that this strategy induces a family of executions \mathcal{E} , each of which is defined by the set of coin flips made by the processors. We can assume that there exists a time τ after which in all executions in \mathcal{E} with non-zero probability no processors send any more messages. Otherwise, the adversary can always wait for another message that must be sent, then for the next message, and so on, until $\Omega(kn)$ messages.

Let us prove that in each execution $E \in \mathcal{E}$ every processor in the bubble must eventually leave the bubble before returning, which implies $\Omega(kn)$ messages in executions in which all processors return. Towards this goal, we first show that a processor cannot return while still in the bubble. Then we prove that all processors in the bubble are forced to either return while still in the bubble (which cannot happen) or leave the bubble, completing the proof.

For the first part, assume for the sake of contradiction that there exists an execution $E \in \mathcal{E}$ and a processor $p \in S$ that decides in E while still being in the bubble. Practically, this implies that p has returned from its method invocation without receiving any messages, and without any of its messages being received. To obtain a contradiction, we build two alternate executions E' and E'' , both of which are indistinguishable to p , but in which p must return different outputs.

In execution E' , we run all processors outside the bubble until one of them returns—this must eventually occur with constant probability, since this execution is indistinguishable to these processors from an execution in which all (at most $k/4 < n/2$) processors in the bubble are initially crashed. We suspend messages sent to the processors inside the bubble. We then run processor p , which flips the same coins as in E (the execution exists as this happens with probability > 0), observes the same emptiness and therefore eventually returns with constant probability, without having received any messages. We deliver p 's messages and suspended messages as soon as p decides.

In execution E'' , we first run p in isolation, suspending its messages. With probability > 0 , p flips the same coins as in E , and must eventually decide with constant probability without having received any messages. We then run all processors outside the bubble in lock-step. One of these processors must eventually return with constant probability, since to these processors, the execution is indistinguishable from an execution in which p (and other processors in the bubble) has crashed initially. We deliver p 's messages after this decision. Since both E' and E'' are indistinguishable to p , it has to return the same value in both executions with constant probability. However, this cannot be the case because instances of method M are

strongly non-commutative, the two returning instances are not concurrent, and occur in opposite orders in the two executions. This correctness requirement is enforced *deterministically*. Therefore, p must return distinct values in executions E' and E'' , which is a contradiction. Hence, p cannot return in E .

To complete the argument, we prove that p has to eventually return or leave the bubble, with probability $\geq \alpha$. We cannot directly require this of the execution prefix E since not all messages by correct processors have been delivered in this prefix. For this, we consider time τ , at which we crash all recipients of messages by p , and all processors that sent messages to p in E . By the definition of the bubble, the number of crashes we need to expend is $< n/4$. Therefore, by definition of τ , there exists a valid execution, in which no more messages will be sent and p must eventually decide with probability $\geq \alpha$. From p 's prospective, the current execution in the bubble can be this execution, and if the adversary keeps p in the bubble for long enough, it has to decide with probability $\geq \alpha$. However, from the previous argument, we know that p cannot decide while in the bubble, therefore p has to eventually leave the bubble in order to be able to decide and return.

This shows that a specific processor p must eventually leave the bubble. The final difficulty is in showing that we can apply the same argument to *all* processors in the bubble at the same time without exceeding the failure budget. Notice however that we could apply the following strategy: for each processor q in the bubble, we could fail all senders and recipients of q ($< n/4$), and also all other processors in the bubble ($< n/4$) at time τ . This can be applied without exceeding the failure budget. Since any processor q could be the sole survivor from the bubble to which we have applied the buffering strategy, and since q does not see a difference from an execution in which it has to return, analogously to the previous case, we obtain that each q in the bubble has to eventually leave the bubble with probability $\geq \alpha$.

Therefore, we obtain that at least $\alpha kn/16$ messages have to be exchanged during the execution, which implies the claim. \square

It is easy to check that the *elect* procedure of a leader election algorithm and the *rename* procedure of a strong renaming algorithm are both non-commutative. (In the case of renaming, consider $n + 1$ distinct processors executing the rename procedure. By the pigeonhole principle, there exists some non-zero probability that two processors choose the same name in solo executions. Therefore, these two operations do not commute, and therefore the *rename* procedure is strongly non-commutative.) We therefore obtain the following corollary.

Corollary B.3. *Any implementation of leader election or renaming by $k \leq n$ processors which ensures termination with probability at least $\alpha > 0$ in an asynchronous message-passing system where $t < n/2$ processors may fail by crashing must have worst-case expected message complexity $\Omega(\alpha kn)$.*