

Early Deciding Synchronous Renaming in $O(\log f)$ Rounds or Less

Dan Alistarh, Hagit Attiya, Rachid Guerraoui, Corentin Travers

► **To cite this version:**

Dan Alistarh, Hagit Attiya, Rachid Guerraoui, Corentin Travers. Early Deciding Synchronous Renaming in $O(\log f)$ Rounds or Less. [Research Report] 2012. <hal-00687555v2>

HAL Id: hal-00687555

<https://hal.inria.fr/hal-00687555v2>

Submitted on 14 May 2012

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Early Deciding Synchronous Renaming in $O(\log f)$ Rounds or Less

Dan Alistarh	Hagit Attiya
EPFL	Technion
<code>alistarh@epfl.ch</code>	<code>hagit@cs.technion.ac.il</code>

Rachid Guerraoui	Corentin Travers*
EPFL	Univ. Bordeaux
<code>guerraoui@epfl.ch</code>	<code>travers@labri.fr</code>

Abstract

Renaming is a fundamental problem in distributed computing, in which a set of n processes need to pick unique names from a namespace of limited size. In this paper, we present the first *early-deciding* upper bounds for synchronous renaming, in which the running time adapts to the actual number of failures f in the execution. We show that, surprisingly, renaming can be solved in *constant* time if the number of failures f is limited to $O(\sqrt{n})$, while for general $f \leq n - 1$ renaming can always be solved in $O(\log f)$ communication rounds. In the wait-free case, i.e. for $f = n - 1$, our upper bounds match the $\Omega(\log n)$ lower bound of Chaudhuri et al. [12].

1 Introduction

Unique names, or identifiers, are a fundamental prerequisite for solving a variety of problems in distributed computing. While in many settings unique names are available, they often come from a very large, practically unbounded namespace, which reduces their usefulness. The *renaming* problem [3], in which a set of processes need to be assigned names from a namespace of small size, is one of the fundamental problems in distributed computing, and a significant amount of work, e.g. [1, 2, 3, 5, 6, 7, 8, 9, 16, 17], studied its solvability and complexity in fault-prone distributed systems.

Much of the work on the renaming problem, whether in shared-memory [1, 5] or in message-passing systems [3], has focused on the asynchronous case, in which the processes' steps or messages may be delayed arbitrarily by the scheduler. However, real world systems experience long periods where delays are bounded, and inter-process communication is synchronous, even though processes may still fail by crashing. The complexity of renaming in such a synchronous setting, where processes communicate by message-passing, was investigated by Chaudhuri et al. [12] and Okun [18]. In brief, they found that for n processes, up to $n - 1$ of which may fail by crashing, there exist algorithms that assign a tight namespace of names $1, \dots, n$ names in $O(\log n)$ rounds of communication. Chaudhuri et al. [12] also showed a matching lower bound.

The analysis in these papers only focused on the case where the *maximum* number of processes, $n - 1$, may fail by crashing during an execution. However, the $\Omega(\log n)$ lower bound of [12] does

*Additional supports from ANR projects ALADDIN and DISPLEXITY

not exclude algorithms that terminate faster when the *actual* number of failures f in the execution is less than $n - 1$. Such speculative algorithms, also known as *early deciding*, are known to exist for *consensus* [19] and *set agreement* [11]. Early-deciding protocols achieve consensus in $O(f)$ communication rounds [14], and k -set agreement in $O(\frac{f}{k})$ rounds [15].

It is therefore natural to ask what is the time complexity of early deciding synchronous renaming in a message-passing system. The answer to this question does not appear trivial: for example, a successful strategy for renaming [12] was for each process to obtain a new bit from its name by running a simple one round symmetry-breaking protocol—after $\Theta(\log n)$ rounds, each process has a unique name from 1 to n . However, it is hard to speed up this approach to obtain more bits for the process’s name in rounds where there are few failures, without breaking name uniqueness. Another approach [3], where each process proposes a name in each round based on what every other process proposed in previous rounds, until there are no collisions, turns out to be very difficult to analyze when the adversary has a limited budget of f failures.

In this paper, we overcome these challenges, and present the first early-deciding upper bounds for renaming. In short, we find that the complexity of the problem is strongly coupled with the relation between n , the number of processes and f , the failure budget of the adversary. We show that there exists an algorithm that ensures a *tight* namespace of n names, and terminates in a *constant* number of rounds in every execution where $f \in O(\sqrt{n})$, and, in $O(\log f)$ rounds, otherwise. The existence of a constant-time renaming algorithm for non-trivial f is surprising, since the early-deciding bounds for consensus [14] and set agreement [15] are *linear* in f irrespective of the relation with the total number of processes n .

Our second result is an algorithm that improves on the constants in the asymptotic notation, terminating in $\log f + 5$ rounds, and assigning names from 1 to $2n$.

The first algorithm is a slight modification of a result by Okun [18]. His protocol is based on a novel connection between synchronous renaming and the approximate agreement problem [13]. In brief, processes assign temporary ranks to each process identifier that they receive, and perform approximate agreement to converge on an approximate rank for each initial process identifier, within a carefully chosen approximation factor. Each process returns this approximate rank, rounded to the nearest integer, as its name: the protocol ensures that, upon termination, the approximate ranks are far enough apart so that no two processes decide on the same rank. Okun [18] showed that this protocol ensures a tight namespace, preserves the order of the initial identifiers, and terminates in $O(\log n)$ rounds in all executions.

Our main contribution is analyzing this protocol for general f , and showing that it terminates in *constant* time if $f \in O(\sqrt{n})$, and in $O(\log f)$ rounds otherwise. Our analysis characterizes the optimal adversarial strategy for arbitrary values of the parameters n and f —we achieve this by carefully bounding the approximation factor of the approximate agreement protocols relative to the number of failures that the adversary expends in each round, showing that this factor goes down fast, and the algorithm terminates very quickly, if the adversary does not fail a significant fraction of the processes in each round.

The second renaming algorithm we present is simpler, and achieves better constants in the asymptotic notation, while relaxing the size of the namespace to $2n$. The protocol is split in two phases: In the first phase, each process identifies a set of names it is interested in, whose size is proportional to the number of failures that the adversary expends in the phase. In the second phase, processes proceed to progressively halve the size of the set of names they are interested in, until each set is a singleton. At this point, each process may adopt the single name in its set. The

key technical difficulty is in assigning each process the “right” set of names at the end of the first phase, so that there are always enough names for the set of participants. To ensure this, we need to relax the total size of the namespace that processes are interested in to be of size $2n$. The halving procedure in the second phase is similar to the $O(\log n)$ round algorithm by Chaudhuri et al. [12].

We therefore show that the time complexity of early-deciding renaming is upper bounded by $O(\log f)$ synchronous rounds for general f , and can be constant for non-trivial $f \in O(\sqrt{n})$. Both algorithms are *adaptive*, since they do not know the number of participants n in advance; they also adapt to the number of failures f in the execution.

Roadmap. We present the problem statement and the model in Section 2, and give an overview of related work in Section 3. In Section 4, we present the analysis of the tight renaming algorithm, while the second algorithm can be found in Section 5. We conclude with a discussion of open questions in Section 6.

2 Model and Problem Statement

Model. We consider a standard synchronous message passing system with n processes p_1, \dots, p_n . Initially, processes have unique identifiers from a namespace of unbounded size. Time is divided into *rounds*, and the processes’ clocks are synchronized. Each round proceeds as a sequence of *send*, *receive*, and *process* steps, in which processes may send and receive messages, and perform local computation if necessary. We assume that at most $t < n$ processes may fail by crashing. If crashed, a process stops taking further steps in the execution; a process may fail to send any subset of its messages in rounds in which it crashes.

Let f denote the number of failures in the current execution. We focus on *early-deciding* algorithms, that adapt their time complexity to the actual number of failures in the execution, i.e. whose running time is a function of f only.

Renaming. The *renaming* problem [3] requires that each correct process eventually returns a name, and that the names returned should be unique. The size of the resulting namespace should only depend on the parameters n . The *tight* renaming problem requires that the size of the namespace be exactly n ; otherwise, we say that the solution is *loose*.

Approximate Agreement. Consider a small real number $\epsilon > 0$. In the ϵ -*approximate agreement* problem [13], each process p_i starts with a proposal v_i , which is a real number. An approximate agreement algorithm must satisfy the following properties: (1) Each correct process eventually decides a value d_i ; (2) For any correct processes p_i and p_j , $|d_i - d_j| \leq \epsilon$; (3) For any correct process p_ℓ , there are processes p_i and p_j with initial values v_i and v_j , such that $v_i \leq d_\ell \leq v_j$.

Notation. Throughout this paper, \log denotes the logarithm base two. Moreover, we write $\log f$ instead of $\lfloor \log f \rfloor$. To simplify notation, we assume that $f \geq 1$, which allows us to consider running times of $O(\log f)$ rounds.

3 Related Work

The renaming problem was introduced in [3], where the authors also provide a wait-free solution using $(2n - 1)$ names in an asynchronous message-passing system, and show that at least $(n + 1)$ names are required in the wait-free case. This lower bound on the namespace size for the case of

asynchronous solutions was improved to $(2n - 2)$ by Herlihy and Shavit [16], and Rajsbaum and Castañeda [10]. (This lower bound holds when n is a power of a prime number.)

Considerable research has analyzed the upper and lower bounds for renaming in asynchronous setting, in particular in shared memory e.g., [1, 2, 5, 6, 7, 8, 17]. For a detailed description of these results, we refer the reader to e.g., [1].

On the other hand, there has been relatively little work on the complexity of synchronous renaming. Herlihy et al. [12] considered wait-free renaming in a synchronous message-passing system, identical to the one we consider in this paper. They prove a lower bound of $\Omega(\log n)$ rounds on the time complexity of the problem in runs where $t = n - 1$ processes may crash, and provide a matching algorithm that achieves tight renaming with time complexity $\lceil \log n \rceil + 3$ in every execution. In contrast, our algorithms are *early deciding*, in that they adapt to the *actual* number of failures f in the current execution, deciding in $O(\log f)$ rounds. On the other hand, the lower bound argument of Herlihy et al. [12] applies to our algorithms as well¹, implying a lower bound of $\Omega(\log n)$ rounds in executions where $f = n - 1$.

In [18], Okun presents a synchronous message-passing algorithm for tight renaming algorithm with $O(\log n)$ time complexity. His algorithm is based on a novel connection between renaming and approximate agreement. In brief, processes perform approximate agreement on the rank of each initial identifier, until they are certain that no two processes obtain the same rank. In this paper, we extend this technique by providing a new analysis for minor variation of his algorithm, proving that its running time is in fact $O(\log f)$ in executions with f failures, which is asymptotically optimal.

4 A Tight Renaming Algorithm

In this section, we analyze the tight renaming algorithm of Okun [18] and prove that it terminates in $O(\log f)$ rounds, where $f < n$ is the number of processes that the adversary crashes in the current execution. We begin with a short description of the algorithm; a detailed exposition can be found in the original paper [18].

Algorithm Overview. The algorithm is based on a novel connection between renaming and the approximate agreement problem. First, a simple synchronous approximate agreement (AA) algorithm is introduced. Then, the algorithm runs in parallel at most n separate instances of this approximate agreement algorithm, one for *each* process in the system. The goal is to agree on an approximate rank for each process's initial identifier, which will be the value decided by the corresponding approximate agreement algorithm, rounded to the nearest integer. The key ingredient is to run the approximate agreement algorithm for long enough to ensure that the decision values of the AA protocols corresponding to each initial identifier are sufficiently spaced, ensuring that the rank decided for each process is unique.

The Approximate Agreement Algorithm. The algorithm, whose pseudocode appears in Figure 1, proceeds as follows. Each process starts with an initial value v_i and a desired approximation factor ϵ , which bounds the maximum desired skew between decided values. The process maintains an estimate d_i of its decision value, which is updated in every round to the arithmetic average of all values received. Once all the estimates received in a round are within at most ϵ of each other, the process returns its current estimate.

¹While the original lower bound of Herlihy et al. [12] applies only to comparison-based algorithms, an argument by Attiya et al. [4] generalizes this bound to a wider class of algorithms, which includes the ones in this paper.

```

1 procedure proposei(vi, ε) ;
2   di ← vi;
3   for each round r ≥ 0 do
4     broadcast(di);
5     δ ← maxℓ≠j(|dℓ - dj|);
6     if δ ≤ ε then return di;
7     else di ← arithmetic average of all values dj received;

```

Figure 1: The Approximate Agreement algorithm.

```

1 procedure renamei(vi) ;
   /* Phase one */
2   for each round r = 1, 2, 3 do
3     broadcast(vi);
4     Cr ← the number of distinct identifiers received in round r;
5   n ← C1;
6   V ← the set of identifiers received in round 3;
7   ε ←  $\frac{1}{10C_2}$ ;
   /* Phase two */
8   for each round r ≥ 4 do
9     for every identifier id ∈ V do
10      Participate in id's instance of the approximate agreement algorithm,
11      with initial value C2 · rankV(id) , until the algorithm returns;
      /* rankV(id) is the rank of id in the set V, in increasing order */
      /* Upon completion of all the approximate agreement algorithms */
12   namei ← final value in vi's instance of the approximate agreement algorithm, rounded to the
      nearest integer;
13   return namei;

```

Figure 2: The Tight Renaming Algorithm.

The Renaming Algorithm. Each process p_i starts with an initial name v_i . The algorithm, whose pseudocode appears in Figure 2, has two phases.

The first phase contains the first three rounds, in which processes exchange their identifiers, in order to identify the parameter n and the relative ranks of their identifiers.

In the second phase, which starts at round four, based on the information computed so far, each process proposes a rank for each participating process to a separate instance of the approximate agreement algorithm in Figure 1. Notice that all these agreement instances run in parallel; all messages by a process in a round (one for each AA protocol in which it participates) are packaged into a single composite message, which each process sends in each round. The approximation factor for all these agreement instances is $1/(10C_2)$, where C_2 is the number of distinct identifiers the process received in round 2 of the first phase. (This factor is chosen such that no two identifiers may receive the same final rank from the approximate agreement instances when rounded to the nearest integer.) The algorithm terminates when all the approximate agreement algorithms terminate, ensuring the desired approximation factor.

Name uniqueness. We give a brief overview of the mechanism ensuring name uniqueness; a complete analysis can be found in [18]. Recall that, for each initial identifier id , each process p proposes the rank of id multiplied by C_2/C_3 as the initial value in id 's instance of the AA protocol. Obviously, the processes' ranks for the same identifier may be distinct (as a consequence of failures in the first phase). However, a key observation is that they will be distinct in a consistent way: if p proposes rank 6 for id α , and rank 7 for id $\beta > \alpha$, then another process q proposing rank 7 for id α will have to propose rank at least 8 for id β . Analyzing the AA protocol, this will imply that, given any process p , its decision values for distinct ids $\alpha < \beta$ will be at distance at least 1 from each other [18, Theorem 3].

This mechanism might still allow the possibility that two distinct processes decide on the same name when rounding their AA decision value to the nearest integer. This is handled by multiplying the initial ranks with C_2/C_3 . This ratio is higher than 1 only when a processor observes crashes between the second and third rounds, and the algorithm ensures that for any ids α and β , their corresponding decision values at any two processes p and q are at distance > 1 from each other [18, Lemma 3]. In turn, this implies that no two processes may decide the same name.

Analysis. Our key observation is that the variant of the synchronous approximate agreement algorithm of Okun [18] presented in Figure 1 guarantees the required approximation factor of $1/(10C_2)$ in a constant number of rounds when $f > \sqrt{n}/2$, and $O(\log f)$ rounds, otherwise. In turn, this implies that the renaming algorithm terminates within three additional rounds.

To prove this, we first introduce some notation. Let $\sigma(S)$ be the *diameter* of a set S , i.e. $\max_{a,b \in S}(|a - b|)$.

For any round $r > 0$ in the execution of the approximate agreement algorithm, let U_r be the multiset of distinct values that processes that are active (i.e., send at least one message) in round r have in the beginning of the round.

Let f_r be the number of processes that crash in round r ; for convenience, denote $f_0 = 0$. Denote by δ_r the fraction of processes that crash in round r from among the processes that did not crash before round r ; that is, $\delta_r = f_r / (n - \sum_{i=0}^{r-1} f_i)$

We first state the following lemma, bounding the diameter of the set U_{r+1} depending on the diameter of U_r and the fraction of processes that crash in round r ; its proof follows [18, Lemma 1].

Lemma 1. $\sigma(U_{r+1}) \leq \frac{2\delta_{r+1}}{1-\delta_{r+1}} \cdot \sigma(U_r)$.

The next lemma bounds the number of rounds needed for the approximate agreement algorithm of Figure 1 to achieve a maximum diameter of $\epsilon = 1/(10n)$ for the set of decisions corresponding to each initial value, when starting with proposal sets of diameter $\leq n$. The bound depends on f , the number of failures that the adversary expends in total, and on the number of failures f_i which the adversary expends in each round $i \geq 1$.

Lemma 2. *Consider an execution of the approximate agreement algorithm of Figure 1, starting with an initial set of diameter $\sigma(U_1) \leq n$, in which at most f processes crash. Let R be the number of rounds needed for the algorithm to reach a diameter $\epsilon \leq 1/(10n)$. The following claims hold:*

- If $f \leq \sqrt{n}/2$, then R is a constant.
- If $\sqrt{n}/2 < f \leq n - 1$, then $R \leq 5 \log f + 10$.

Proof. We assume $n \geq 6$; the claim can be checked for $n \leq 5$ by calculation.

The first case is when $f \leq \sqrt{n}/2$. In this case, notice that the fraction of processes that fail in any round of the protocol is at most \sqrt{n}/n , i.e. $\delta_r \leq \sqrt{n}/n$, for all $r \geq 1$. In turn, by Lemma 1, the diameter of the set of values for the current instance of approximate agreement is reduced by at least $2\delta_r/(1-\delta_r) \leq 1/(\sqrt{n}-0.5)$ in each round r . Therefore, after the first 10 rounds, the diameter of this set is at most

$$\frac{n}{(\sqrt{n}-0.5)^{10}} \leq \frac{1}{10n}, \text{ for any } n \geq 6.$$

Therefore the maximum diameter the end of round 10 is $\leq 1/(10n)$, as claimed.

In the second case, we assume that $f > \sqrt{n}/2$. First, notice that δ_r , the fraction of active processes that crash in a round r , can be greater than $1/2$ in at most $\lfloor \log_2(f+1) \rfloor + 1$ distinct rounds. Therefore, any execution of at least $5\lfloor \log_2(f+1) \rfloor + 10$ rounds contains at least $4\lfloor \log_2(f+1) \rfloor + 9$ rounds r in which $\delta_r < 1/2$. Lemma 1 implies that the diameter of the set U_r at the end of these rounds is at most

$$\frac{n}{2^{4\lfloor \log_2(f+1) \rfloor + 9}} \leq \frac{n}{2^9 \cdot (f+1)^4} \leq \frac{1}{10n} \text{ for all } n \geq 1,$$

where the last step uses the fact that $f \geq \sqrt{n}/2$. Therefore, in this case, the number of rounds necessary to obtain a maximum diameter of at most $1/(10n)$ is $O(\log f)$. \square

We conclude that the resulting renaming algorithm is early deciding, terminating in a constant number of rounds, when $f \leq \sqrt{n}/2$, and $O(\log f)$ rounds, otherwise.

Theorem 3. *For any f , $0 \leq f \leq n-1$, let $\ell_f = \text{constant}$ if $f \leq \sqrt{n}/2$, and $\ell_f = 5\log(f+1) + 10$, otherwise. Then the renaming algorithm in Figure 2 is a tight order-preserving renaming algorithm with time complexity $O(\ell_f)$ and message complexity $O(n^2\ell_f)$ in executions where f processes fail by crashing.*

Proof. We focus on early decision, since the other properties follow from Theorem 4 of [18]. First, recall that the initial value v for each approximate agreement algorithm is computed locally by each process as $\text{rank}_V(\alpha) \cdot C_2/C_3$. Since $\text{rank}_V(\alpha)/C_3 \leq 1$, by the definition of C_2 , we obtain that all initial values are between 1 and n . This also implies that our desired approximation factor ϵ is at most $1/(10n)$.

Lemma 2 implies that all instances of approximate agreement that the renaming algorithm executes in parallel terminate in $O(\ell_f)$ rounds, with an approximation factor $\epsilon \leq 1/(10n)$. We obtain that the renaming algorithm terminates in $O(\ell_f)$ rounds, as claimed. \square

5 A Loose Renaming Algorithm with Improved Round Complexity

In this section, we present a loose renaming algorithm that terminates in $\log f_1 + 5$ rounds and uses at most $2n$ names, where n is the number of participating processes and f_1 the number of failures that occur in the first round. Our algorithm extends the non-early deciding renaming algorithm by Chaudhuri, Herlihy, and Tuttle [12]. The latter algorithm is tight; its round complexity, however, depends solely on n and not on the number of failures among the participating processes. Specifically, the namespace when n processes participate is $[1, n]$ and the algorithm terminates in $\log n + 2$ rounds. In the following, the algorithm of [12] is called the *CHT-renaming* algorithm.

Algorithm Overview. The pseudo-code of the algorithm appears in Figure 3. It contains two phases.

In the first phase, which consists of the two first rounds (line 1-line 7), each process selects an interval of names in which it wishes to pick its final name. The size of each interval is upper-bounded by $4f_1$, where f_1 is the number of processes that fail during the first round. The second phase (line 8-line 16) consists of a variant of the CHT-renaming algorithm. Processes use this procedure to progressively shrink the interval of names they are interested in, until obtaining an interval of size 1. The algorithm ensures that no two processes obtain the same interval of size 1, i.e. a single name. Thus, each process can decide the unique name in its final interval. Moreover, it guarantees that in each round, processes holding the largest intervals reduce their interval by one half. Therefore our algorithm terminates in $O(\log f_1)$ rounds.

We begin with a brief description of the CHT-renaming algorithm and then explain how each process selects an initial interval of names.

Definitions and notations. Before describing the algorithm in more details, we introduce some notations, extending those in [12]. An interval I of positive integers is *well-formed* if I is of the form $[d2^j + 1, (d + 1)2^j]$ for some positive integers d, j . Note that for every pair I, J of well-formed intervals, $I \cap J \neq \emptyset \implies I \subseteq J \vee J \subseteq I$. Given a set \mathcal{I} of well-formed intervals, we say that interval $I \in \mathcal{I}$ is *maximal* in \mathcal{I} if for every $J \in \mathcal{I}$, either $I \cap J = \emptyset$ or $J \subseteq I$.

The CHT-renaming Algorithm Each process p maintains a well-formed interval of names I , which are the names p is interested in. In each round, process p sends its id and the interval it currently holds I . The intervals intersecting with I that p receives are stored in the set \mathcal{I} (line 10). p also stores the set of ids of the processes that are conflicting with p , that is, processes that hold an interval intersecting with I in the set \mathcal{P} (line 11). If I is maximal in \mathcal{I} (line 12), I is split in half, and p picks the bottom half or the top half of I , denoted $bot(I)$ and $top(I)$ respectively, as its new interval (line 13-line 14). More precisely, the new interval of p is $bot(I)$ if the rank of p 's identity in \mathcal{P} is smaller than $\frac{|I|}{2}$. Otherwise, p selects $top(I)$ as its new interval. Finally, if p observes that every interval sent in the round has size 1, it decides the unique name in its interval. Since in each round the size of maximal interval is at least divided by 2, the algorithm terminates after $O(\log c)$ rounds, where c is the size of the largest initial maximal interval.

Name uniqueness relies on the following invariant, which is satisfied by the set of intervals $\mathcal{I}[r]$ held by the processes at the beginning of each round r :

Invariant 1. *For every $I \in \mathcal{I}[r]$, if m processes hold an interval $I' \subseteq I$, then $m \leq |I|$. In particular, this means that I is large enough to allow processes with an interval $I' \subseteq I$ to decide distinct names in the interval I .*

In the original CHT-renaming, the initial interval of each process p is of the form $[1, 2^b]$, where 2^b is the least power of 2 larger than or equal to the number of participating processes from which p has received a message in the first round. The invariant above is thus initially true. In our algorithm, initial intervals are selected differently, but the invariant is still satisfied.

Selection of initial intervals. In the first round, process p broadcasts its id, and ranks its id among the ids it receives (line 2). Let rk_p be the rank obtained by p . If, among the ids of the participating processes, the rank of the id of p is i , and f_1 is the number of failures that occur in the first round, then

$$i - f_1 \leq rk_p \leq i. \tag{1}$$

```

1 procedure renamei(idi) ;
   /* Round 1 */
2 broadcast(idi);
3 rk ← the rank of id among the id received in round 1;
   /* Round 2 */
4 broadcast(idi, rk);
5 let rk_max = largest rank received; for each j, 1 ≤ j ≤ rk_max do hj ← |{(id', rk') : rk' = j}| ;
6 est_f ← max{(∑j ∈ I hj) - |I| : I ⊆ [1, rk_max]};
7 let j = ⌈log(est_f)⌉ and d such that d · 2j + 1 ≤ rk ≤ (d + 1) · 2j ;
8 I ← [d · 2j + 1, (d + 1) · 2j];
   /* Round 3, 4, ... : CHT-renaming protocol [12] */
9 repeat
10 broadcast(id, I) ;
11 I ← {I' : ⟨id', I'⟩ received and I ∩ I' ≠ ∅} ;
12 P ← {id' : ⟨id', I'⟩ received and I ∩ I' ≠ ∅} ;
13 if ∀I' ∈ I, I' ⊆ I then
14 let bot(I) and top(I) the bottom half and top half of I respectively;
15 if rank(id, P) ≤  $\frac{|I|}{2}$  then I ← bot(I) else I ← top(I)
16 until ∀I' ∈ I, |I'| = 1 ;
17 return namei where I = [namei, namei] ;

```

Figure 3: The Loose Renaming Protocol.

This is because p may miss at most f_1 messages from processes with id smaller than i . In the second round, p sends its rank together with its id. It then estimates, based on the ranks it receives, the number of failures that occur in the first round. To that end, p evaluates for each interval of names I the difference between the number of processes ranked in I and the size of I (line 4-line 5). The estimate est_f of the number of failures is then the maximum over all differences. More precisely,

$$est_f = \max_{I \subseteq [1, rk_max]} \sum_{i \in I} h_i - |I|.$$

where rk_max is the largest rank received by p and h_i is—to the knowledge of p —the number of processes that rank their id i in the first round. The estimation est_f is upper-bounded by f_1 , the number f_1 of failures that occur in the first round. To see why, consider an interval $I = [a, b]$. We have:

$$\sum_{i \in I} h_i \leq |\{q : a \leq rk_q \leq b\}| \leq |\{q : a - f_1 \leq rank(id_q, P) \leq b\}| \leq |I| + f_1, \quad (2)$$

where P is the set of participating processes and $rank(id_q, P)$ is the rank of id_q among the ids of the processes in P . The last inequality follows from Equation 1.

Finally, p selects a well-formed interval. This is performed in two steps. First, p chooses two integers d, j such that the well-formed interval $[d2^j + 1, (d + 1)2^j]$ contains rk_p , and 2^j is the least power of 2 larger than or equal to est_f . However, by equation (2), at most $2^j + est_f \leq 2^{j+1}$ processes may have their rank contained in $[d2^j + 1, (d + 1)2^j]$. Then, in order to satisfy Invariant 1, p selects the interval $I = [d2^{j+1} + 1, (d + 1)2^{j+1}]$ as its initial interval. Since a process p' that selects an interval $\subseteq I$ has its rank $rk_{p'}$ contained in $[d2^j + 1, (d + 1)2^j]$, at most $2^{j+1} = |I|$ processes select

intervals $I' \subseteq I$. Invariant 1 is thus satisfied, which preserves the correctness of the CHT-renaming algorithm. Finally, notice that the size of each initial interval is at most $4f_1$, where f_1 is the number of failures in the first round. Hence the total running time is $O(\log f_1)$ rounds.

Proof of correctness. We first establish some properties of the intervals selected in round 2 (Lemma 4 and Lemma 5). The first lemma shows that the intervals selected in round 2 are “large enough” in the following sense: if I is selected, then at most $|I|$ processes select an interval $I' \subseteq I$. This lemma is the key to showing name uniqueness, i.e., no two processes decide on the same name.

Lemma 4. *Suppose that $\langle id, I \rangle$ is selected in round 2. Let P_I be the set of processes that select an interval $I' \subseteq I$ in round 2. Then $|P_I| \leq |I|$.*

Proof. Let p denote a process selects interval I in round 2 (line 7). By the code, I is thus of the form $[d2^{j+1}+1, (d+1)2^{j+1}]$. Consider a process p' that sends an interval $I' = [d'2^{j'+1}+1, (d'+1)2^{j'+1}] \subseteq I$ in round 3. Denote by rk' the rank obtained by p' in the first round (line 2). By definition of d' (line 6), $d'2^{j'} + 1 \leq rk' \leq (d' + 1)2^{j'}$, and as $I' \subseteq I$, it follows that $d2^j + 1 \leq rk' \leq (d + 1)2^j$.

Notice also that, as p' selects I' in round 2, a message $\langle id', rk' \rangle$ from p' must have been received by p in round 2. Therefore,

$$|P_I| \leq |\{\langle id', rk' \rangle : \langle id', rk' \rangle \text{ is received by } p \text{ and } d2^j + 1 \leq rk' \leq (d + 1)2^j\}|$$

Recall that, at process p , h_{rank} denotes the number of messages of the form $\langle id, rank \rangle$ received by p in round 2 (line 4). It thus follows from the definition of the number of failures est_f estimated by p (line 5) that:

$$|P_I| \leq \sum_{d2^{j+1}+1 \leq i \leq (d+1)2^{j+1}} h_i \leq (d + 1)2^j - d2^j + est_f \leq 2^{j+1}$$

The last inequality follows from the fact that, at process p , $j = \lceil \log est_f \rceil$. As $I = [d2^{j+1} + 1, (d + 1)2^{j+1}]$, $|I| = 2^{j+1}$. Therefore $|P_I| \leq |I|$. \square

The next lemma is the basis to establish that the size of the namespace is at most twice the number of participating processes. It is the analogue of Lemma 5 in [12].

Lemma 5. *If $I = [a, b]$ is a maximal interval selected by some process round 2, then there are least $\frac{a-1}{2}$ processes that either select an interval $[a', b']$ with $b' < a$ at the beginning of round 2 or have failed by the end of round 2.*

Proof. Let p denote a process that selects I in round 2. Since I is selected in round 2, $a = d2^{j+1} + 1$ for some $d, j \geq 0$. Moreover, the rank rk_p of the identity id_p of p according to the ids p receives in round 1 satisfies $d2^j + 1 \leq rk_p$ (line 6). Hence, p receives at least $d2^j = \frac{a-1}{2}$ identities strictly smaller than or equal to $d2^j$ in round 1. Denote by $Q = \{q_1, \dots, q_m\}$ the set of processes with id $\leq \frac{a-1}{2}$ from which p receives a message in round 1. Note that $m \geq \frac{a-1}{2}$. For each process $q_i \in Q$ that has not failed by the end of round 2, let rk_i the rank of id_i among the identities received by q_i in round 2. We consider two cases:

- The rank of each process $q_i \in Q$ is smaller than or equal to $\frac{a-1}{2}$. That is, for each $q_i \in Q$, $rk_i \leq \frac{a-1}{2} = d2^j$ or q_i has failed by the end of round 2. Hence, if $[a_i, b_i]$ is the interval selected by q_i , $a_i < a = 2d^{j+1} + 1$ (line 7). Moreover, as $[a_i, b_i]$ and $[a, b]$ are well-formed intervals, and $[a, b]$ is a maximal interval, $b_i < b$. Hence, in this case, at least $\frac{a-1}{2}$ processes select an interval $[a', b'], b' < a$ or fails by the end of round 2.

- There exists a process $q_i \in Q$ such that $rk_i > \frac{a-1}{2}$. Without loss of generality, assume that $id_1 < \dots < id_m$, where id_i denotes the id of q_i . Let $q_\ell \in Q$ be the process with the smallest identity for which $rk_\ell > \frac{a-1}{2}$. q_ℓ thus receives in round 1 at least $\frac{a-1}{2}$ identities smaller than its own identity id_ℓ . Since p receives $\ell - 1$ ids smaller than id_ℓ , at least $\frac{a-1}{2} - (\ell - 1)$ of those identities are not received by p . Hence, at least $\frac{a-1}{2} - \ell + 1$ processes $\notin Q$ fail in the first round. Moreover, by definition of q_ℓ , each process $q_i, 1 \leq i \leq \ell - 1$ either fails or has rank $rk_i < rk_p$. Hence, by the first item above, if q_i does not fail, the interval $[a_i, b_i]$ it selects is such that $b_i < a$. Therefore, in total $\frac{a-1}{2} - \ell + 1 + (\ell - 1) = \frac{a-1}{2}$ processes fail by the end of round 2 or select an interval $[a', b']$ with $b' < a$. \square

In rounds $r \geq 3$, processes apply the CHT-renaming algorithm using the intervals selected in round 2 as initial intervals. The next two lemmas generalize Lemma 4 and Lemma 5, respectively, to the intervals held by the processes at the beginning of each round $r \geq 3$.

Lemma 6 (Lemma 4 in [12]). *Suppose that $\langle id, I \rangle$ is sent in round $r \geq 3$. Let P_I be the set of processes that send an interval $I' \subseteq I$ in round r . $|P_I| \leq |I|$.*

Proof. For $r = 3$, the lemma immediately follows from Lemma 4, as each interval sent in round 3 is selected in round 2. For $r > 3$, the proof can be found [12]. The proof depends solely on the code of the CHT-renaming algorithm and the fact that the lemma is true in the rounds preceding r . \square

Lemma 7 (Lemma 5 in [12]). *If $I = [a, b]$ is a maximal interval received by some process p during round $r \geq 3$, then there are at least $\frac{a-1}{2}$ processes that either hold an interval $[a', b']$ with $b' < a$ at the beginning of round r or has failed by the end of round r .*

Proof. Again the lemma immediately follows from Lemma 5 for $r = 3$. For $r > 3$, we rely on the proof of a similar lemma (Lemma 5) in [12]. There, it is shown that $I = [a, b]$ is a maximal interval received by some process p during round $r \geq 3$, then there are at least $a - 1$ processes that either hold an interval $[a', b']$ with $b' < a$ at the beginning of round r or has failed by the end of round r . A close examination of the proof reveals the same argument (replacing $a - 1$ by $\frac{a-1}{2}$) can be used to prove the lemma. \square

In each round, processes that hold maximal intervals split their interval in half. As a result, the size of the largest interval is halved in each round. This immediately follows from Lemma 6 in [12].

Lemma 8. *For $r \geq 3$, let ℓ_r denote the size of the largest interval sent in round r . Then $\ell_r \geq 2\ell_{r+1}$.*

We have now all the ingredients to establish that the algorithm described in Figure 3 solves loose renaming in $O(\log f)$ rounds.

Theorem 9. *The algorithm described in Figure 3 terminates in $\log f_1 + 5$ rounds, where f_1 is the number of failures in the first round. Each non-faulty process decides a distinct name in the interval $[1, 2n]$, where n is the number of participating processes.*

Proof. Termination. Consider an interval I selected by some process p in round 2. By the code, $I = [d2^{j+1} + 1, (d+1)2^{j+1}]$ for some $d, j \geq 0$ and we have $j = \lceil \log f \rceil$, where f is the number of failures that occur in the first round, as estimated by p . As we have seen in the description of the

algorithm (paragraph Selection of initial intervals), $est_f \leq f_1$. Therefore, $|I| = 2 \times 2^j \leq 4est_f \leq 4f_1$. Therefore, ℓ_3 , the size of the largest interval sent in round 3 is smaller than or equal to $4f_1$. Then, as by Lemma 8, the length of the largest interval is divided by two in each round, the size of the largest interval sent in round $\log f_1 + 5$ or earlier is 1. As each non-faulty decides if only interval of size 1 are sent in a round, the algorithm terminates in at most $\log f_1 + 5$ rounds.

Name uniqueness. The proof is the same as in [12]. Suppose that two processes p and p' decide the same name a in round r and r' respectively. Suppose without loss of generality that $r \leq r'$. In the round in which they decide, p and p' send the interval $I = [a, a]$, and therefore interval containing $[a, a]$ in every preceding round. Since, when p decides, every interval that contains I received by p is equal to I (line 16), p' sends also $[a, a]$ in round r . Hence two processes send the same interval of size 1 in round r : a contradiction with Lemma 6.

Namespace. Again, the proof closely follows the proof in [12]. Let p be the process that decides the highest name a . In the round $r \geq 3$ in which p decides, it sends the interval $[a, a]$. Therefore, by Lemma 7, at least $\frac{a-1}{2}$ processes other than p either sends an interval $[a', b']$, with $b' < a$ in round r or fail by the end of round r . The number of participating processes n is thus at least $n \geq \frac{a-1}{2} + 1$, from which we conclude that $a \in [1, 2n]$. \square

6 Discussion

This paper presents the first early-deciding upper bounds for synchronous renaming. We show that, surprisingly, renaming can be solved in *constant* number of rounds if the number of failures f is limited to $O(\sqrt{n})$, while in the general case, renaming can always be solved in $O(\log f)$ communication rounds. In the wait-free case, i.e. for $f = n - 1$, this upper bound is matched asymptotically by the $\Omega(\log n)$ lower bound of Chaudhuri et al. [12]. It remains an open question whether this is tight for other values of f .

References

- [1] Dan Alistarh, James Aspnes, Keren Censor-Hillel, Seth Gilbert, and Morteza Zadimoghaddam. Optimal-time adaptive strong renaming, with applications to counting. In *PODC'11: Proceedings of the 30th annual ACM symposium on Principles of distributed computing*, pages 239–248, 2011.
- [2] Dan Alistarh, James Aspnes, Seth Gilbert, and Rachid Guerraoui. The complexity of renaming. In *FOCS'11: Proceedings of the 52nd Annual IEEE Symposium on Foundations of Computer Science*, pages 718–727, 2011.
- [3] Hagit Attiya, Amotz Bar-Noy, Danny Dolev, David Peleg, and Rudiger Reischuk. Renaming in an asynchronous environment. *Journal of the ACM*, 37(3):524–548, 1990.
- [4] Hagit Attiya and Taly Djerassi-Shintel. Time bounds for decision problems in the presence of timing uncertainty and failures. *Journal of Parallel and Distributed Computing*, 61(8):1096–1109, 2001.
- [5] Hagit Attiya and Arie Fouren. Adaptive and efficient algorithms for lattice agreement and renaming. *SIAM Journal on Computing*, 31(2):642–664, 2001.

- [6] Elizabeth Borowsky and Eli Gafni. Immediate atomic snapshots and fast renaming. In *PODC '93: Proceedings of the twelfth annual ACM symposium on Principles of distributed computing*, pages 41–51, New York, NY, USA, 1993. ACM.
- [7] Alex Brodsky, Faith Ellen, and Philipp Woelfel. Fully-adaptive algorithms for long-lived renaming. *Distributed Computing*, 24(2):119–134, 2011.
- [8] James E. Burns and Gary L. Peterson. The ambiguity of choosing. In *PODC '89: Proceedings of the eighth annual ACM Symposium on Principles of distributed computing*, pages 145–157, New York, NY, USA, 1989. ACM.
- [9] Armando Castañeda and Sergio Rajsbaum. New combinatorial topology bounds for renaming: The upper bound. *Journal of the ACM*, 59(1), March 2012.
- [10] Armando Castañeda and Sergio Rajsbaum. New combinatorial topology bounds for renaming: The lower bound. *Distributed Computing*, 22(5-6):287–301, 2010.
- [11] Soma Chaudhuri. More choices allow more faults: Set consensus problems in totally asynchronous systems. *Information and Computation*, 105(1):132–158, 1993.
- [12] Soma Chaudhuri, Maurice Herlihy, and Mark R. Tuttle. Wait-free implementations in message-passing systems. *Theoretical Computer Science*, 220(1):211–245, 1999.
- [13] Danny Dolev, Nancy A. Lynch, Shlomit S. Pinter, Eugene W. Stark, and William E. Weihl. Reaching approximate agreement in the presence of faults. *Journal of the ACM*, 33:499–516, May 1986.
- [14] Danny Dolev, Rudiger Reischuk, and H. Raymond Strong. Early stopping in byzantine agreement. *Journal of the ACM*, 37(4):720–741, October 1990.
- [15] Eli Gafni, Rachid Guerraoui, and Bastian Pochon. The complexity of early deciding set agreement. *SIAM Journal on Computing*, 40:63–78, January 2011.
- [16] Maurice Herlihy and Nir Shavit. The topological structure of asynchronous computability. *Journal of the ACM*, 46(2):858–923, 1999.
- [17] Mark Moir and James H. Anderson. Fast, long-lived renaming (extended abstract). In *WDAG '94: Proceedings of the 8th International Workshop on Distributed Algorithms*, pages 141–155, London, UK, 1994. Springer-Verlag.
- [18] Michael Okun. Strong order-preserving renaming in the synchronous message passing model. *Theoretical Computer Science*, 411(40-42):3787–3794, 2010.
- [19] Marshall Pease, Robert Shostak, and Leslie Lamport. Reaching agreement in the presence of faults. *Journal of the ACM*, 27(2):228–234, April 1980.