

An algorithm for the asynchronous *Write-All* problem based on process collision^{*}

Jan Friso Groote^{1,2}, Wim H. Hesselink³, Sjouke Mauw^{1,2}, Rogier Vermeulen¹

¹ Eindhoven University of Technology, P.O. Box 513, NL-5600 MB Eindhoven, The Netherlands. e-mail: sjouke@win.tue.nl

² CWI, P.O. Box 94079, NL-1090 GB Amsterdam, The Netherlands. e-mail: jfg@cwi.nl

³ University of Groningen, P.O. Box 800, NL-9700 AV Groningen, The Netherlands. e-mail: wim@cs.rug.nl

Summary. The problem of using P processes to write a given value to all positions of a shared array of size N is called the *Write-All* problem. We present and analyze an asynchronous algorithm with work complexity $\mathcal{O}(N \cdot P^{\log(\frac{x+1}{x})})$, where $x = N^{1/\log(P)}$ (assuming $N = x^k$ and $P = 2^k$). Our algorithm is a generalization of the naive two-processor algorithm where the two processes each start at one side of the array and walk towards each other until they collide.

Key words: write-all problem – wait-free – distributed algorithms – work complexity – PRAM – dynamic load balancing.

1 Introduction

The *Write-All* problem is defined as follows. Use P processes (or processors) to write a given value to all positions of a shared array of size N . Without loss of generality, we shall assume that the array is an integer array and that 1 is the value to be written to all its positions.

If the processes are reliable and run equally fast, it is easy to come up with straightforward, optimal solutions for this problem. The situation is quite different, however, if processes can be faulty or run at widely varying speeds while at least one process remains active. Kadem et al. [10] have shown that under these circumstances an $N + \Omega(P \log N)$ lower bound exists on the amount of work processes must carry out when processes can fail. This means that even if all processes run fully in parallel and no process is actually failing, at least $\Omega(\log N)$ time is required to set the array.

The original motivation for the *Write-All* problem comes from [7]. Here it was shown that any program on a

P process synchronous PRAM (Parallel Random Access Machine) running in time T can be executed on any unreliable PRAM with work complexity $W(P) \cdot T$ where $W(P)$ is the work complexity of any algorithm solving the *Write-All* problem of size P . In [9] an overview is given of the algorithms and PRAM simulations that have been developed so far.

Our motivation is quite different. It comes from the design of wait-free or asynchronous algorithms [4–6], to obtain fast, reliable programs for general purpose parallel computers with typically a few dozen processes that run under widely varying loads.

A common problem on such machines is to carry out a task, consisting of N independent subtasks, with P processes, as quickly as possible. Such tasks are, for instance, copying an array, searching an unordered table, and applying a function to all elements of a matrix. We encountered this problem when we had to find a parallel solution to refresh a hashtable by copying all valid elements to a new array [4].

If we abstract from the nature of the subtasks, the problem of executing N independent tasks is adequately characterized by the *Write-All* problem.

In this paper we present a rather straightforward algorithm to solve the *Write-All* problem on an asynchronous PRAM, i.e. a machine on which the processes can be stopped and restarted at will. This means that it is also suitable for all other fault models as mentioned in Kanellakis and Shvartsman, page 13 [9]. Using different terminology we can say that our algorithm is wait-free, which means that each non-faulty process will be able to finish the whole task, within a predetermined amount of steps, independent of the actions (or failures) of other processes.

For a shared array of size N and P processes, our algorithm has to carry out $\mathcal{O}(N P^{\log(\frac{x+1}{x})})$ amount of work where $x = N^{\frac{1}{\log P}}$. The complexity of parallel algorithms is generally characterized by the total amount of steps that all processes must execute, which is called the *work* of the algorithm, instead of the execution time, which under ideal circumstances, can be obtained by dividing the work by the number of available processes. It

^{*} Corresponding author: Dr. S. Mauw, Department of Mathematics and Computing Science, Eindhoven University of Technology, P.O. Box 513, NL-5600 MB Eindhoven, The Netherlands, Tel. +31 40 2472908, (secr. 2474124), Fax +31 40 2475361, E-mail sjouke@win.tue.nl

should be noted that the worst case behaviour leading to the upper bound $\mathcal{O}(N P^{\log(\frac{x+1}{x})})$ can only be achieved under a rare lock step scenario of the processes. So, we expect average complexity to be much better, which has been confirmed by experiment. In order to understand the asymptotic behaviour of the bound $\mathcal{O}(N P^{\log(\frac{x+1}{x})})$ it is interesting to look at the role of x . The value x roughly denotes the ratio between the length of the array (N) and the number of processes (P). In the limit case, where x goes to infinity, the work complexity converges to $\mathcal{O}(N)$. This means that if there are only few processes, relative to the length of the array, complexity of our algorithm becomes linear in the length of the array.

There are a number of existing solutions to the *Write-All* problem (see [9] for an excellent overview). We compare our algorithm to the algorithms X , X' , AW , AW^T and Y that are all suitable for asynchronous PRAMs, ignoring the solutions suitable for more restricted fault models. From certain perspectives our algorithm improves upon all of these.

Algorithm X is the first asynchronous algorithm for the *Write-All* problem [3]. It is designed for the situation where $P \geq N$ and has work $\mathcal{O}(N P^{\log(\frac{3}{2})})$. In [9] a generalisation of X , called X' is presented for the case $P \leq N$ which has the same upper bound $\mathcal{O}(N P^{\log(\frac{3}{2})})$ for the amount of work. For $N = P$ the algorithm presented here has the same upper bound as X' . For $P < N$ our algorithm is an improvement over X' , since for $P < N$ we have $x > 2$ and then $\mathcal{O}(N P^{\log(\frac{x+1}{x})})$ is better than $\mathcal{O}(N P^{\log(\frac{3}{2})})$.

In [1] two particularly clever algorithms are proposed, called AW and AW^T .

Algorithm AW requires work $\mathcal{O}(P^2 + N \log P)$. When $P \leq \sqrt{N}$ this reduces to $\mathcal{O}(N \log N)$ which is particularly good. However, this bound can only be achieved assuming that a set of permutations of $1 \dots P$ with a specific property is given, which requires exponential time to calculate. Such a set can be generated at random, but then the result 'only' holds with high probability. In order to overcome this problem algorithm Y has been proposed [8]. Algorithm Y is conjectured to have (non probabilistic) work upper bound $\mathcal{O}(N \log N)$, which is confirmed by experiments, but which is unproven.

Algorithm AW^T needs work $\mathcal{O}(q N P^\epsilon)$ where $\epsilon = \log_q \log q^c$ for some constant q that can be freely chosen, and a constant c which, according to the proof in [9], can be chosen to be 2. As $\log_q \log q^2$ goes to 0 when q goes to infinity, algorithm AW^T has superior complexity. However, the constant amount of work that must be done in the preprocessing phase (which is independent of N and P) is exponential in q (see [1]). In order to outperform algorithm X' for any N and P , it must be the case that $\epsilon < \log(\frac{3}{2})$. From this it follows that q must be larger than 80. Therefore, to outperform our algorithm, q must be chosen even larger. In the setting for which we developed our algorithm, we generally have $P < \sqrt{N}$ (and thus $x > 4$), so one must choose $\epsilon < \log \frac{5}{4}$ to make algorithm AW^T perform better than our algorithm. This

means that q needs to be larger than 10^5 . This is the reason why we expect that our algorithm performs much better under practical circumstances.

The present paper has the following structure. In Section 2 we present the algorithm. In Section 3 we prove its correctness and show space and time bounds. Section 4 contains some considerations on using a non-uniform tree as the shared data structure. Finally, Section 5 is reserved for conclusions and further considerations.

2 A collision-based algorithm

2.1 Basic case

Although the asynchronous *Write-All* problem in its general setting is far from trivial, the case that there are only two processes ($P = 2$), allows for a very intuitive and optimal solution. This algorithm solves the problem for any value of N in $N + o(N)$ steps. One process starts at the left of the array and walks to the right, in the meanwhile setting the values of the array elements encountered to 1. The other process does the same from right to left. If the two processes collide, the whole array is processed and the processes can stop. In the worst case, one element of the array is processed twice. We call this algorithm the *Basic Collision* algorithm.

In [3] an extension of this algorithm is described, which works with three processes. It is called algorithm T . Two processes have the same behaviour as described above, but the third process behaves differently. It starts in the middle of the array and fills the array alternately to the left and to the right. If the first two processes collide, it means that the whole array is processed. If, e.g., the first and the third process collide, it means that the left part of the array is processed. Therefore they move to the segment of the array that is not processed yet. The first process starts at the left of this segment, the third process starts again in the middle of this segment, and the second process is still busy filling the segment from the right. This procedure repeats until the array is completely processed. This algorithm is also optimal and the work of this algorithm, measured in terms of actual elements processed, is $N + o(N)$. Algorithm T does not appear to be generalizable to larger numbers of processes.

2.2 Generalized case

Our algorithm generalizes the *Basic Collision* algorithm in a different way. We will call it the *Generalized Collision* algorithm. It is best explained by looking at a simple example with four processes ($P = 4$). We choose $N = 25$ in our example.

The processes operate in pairs. Every pair of processes executes the *Basic Collision* algorithm on successive segments of the array. Each segment has length 5, so there are 5 segments. The four processes start at the



Fig. 1. Initial configuration

locations indicated in Figure 1. The arrows indicate the direction in which each process traverses the segment.

Every time that a segment of the array has been processed by a pair, operation continues at the next segment. The first process of a pair to finish a segment can directly continue with the next segment, without having to wait for the other process. In this way, the pairs walk towards each other through the array in steps of length 5 until they collide. A typical path of the four processes in our example is shown in Figure 2. This figure shows just one possible path, in which all processes roughly operate at the same speed. The algorithm, however, is completely robust with respect to process delays, failures and restarts. This is because every process potentially visits all array elements. As long as one process survives, the whole array will be processed.

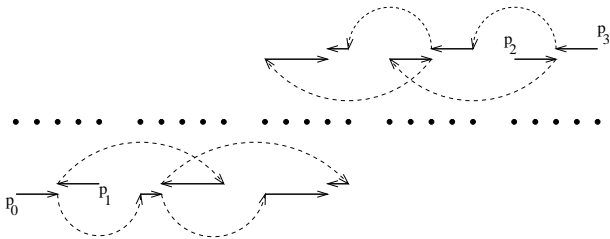


Fig. 2. Possible paths of the processes

From a higher point of view, the four processes also execute the *Basic Collision* algorithm where the grain size of the work is 5. To see this, we have to consider every pair as a single aggregated process and every segment of length 5 as a single aggregated array element. A collision now takes place at a complete segment, rather than at a single array element. This explains why the middle segment in Figure 2 is processed twice.

It is now clear how to generalize this example if we double the numbers of processes and assume 125 array elements. We simply add one level to the hierarchy and have clusters of four processes operate on segments of length 25, until the clusters collide.

This implies that our algorithm works for any number of processes which is a power of two, so $P = 2^k$ for some $k \geq 1$. Furthermore, we have that the length of the array is the length of a basic segment to the same power, so $N = x^k$ for some $x \geq 2$. In the above example we have chosen $k = 2$ and $x = 5$.

In Figure 3 the generalization of the *Basic Collision* algorithm is illustrated in a cube which has to be filled with 1's by 8 processes. The picture shows pairs of processes, clusters of 2 processes, and clusters of 4 processes racing each other. In this example $k = 3$ (the dimension of the cube) so there are 8 processes, and the length of an edge of the cube is x , so that there are x^3 cells to be filled. This is the biggest example that we can easily vi-

sualize in this way. An example with 16 processes would require a 4-dimensional figure.

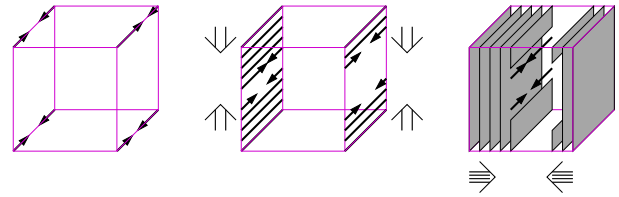


Fig. 3. Generalization of the collision principle illustrated in a cube.

2.3 Data structures

Additional data structures are needed in order to enable the processes to decide which array element should be processed next. First of all, every process has a process identifier (*pid*) consisting of a bit string of length k . The set of all process identifiers is called *PID*. We use the functions *head* and *tail* to return the first element of a bit string and the bit string with the first element deleted. The bit strings will be used to direct the processes to different parts of the array. There is a nice relation between the pids of the processes and the initial position of the processes in the cube from Figure 3. If we consider the general Boolean k -dimensional hypercube, the pids correspond to the processes' initial co-ordinates.

Next, we assume that the processes share a tree of depth k . According to the above explanation, the tree should have a uniform fan-out x . This means that there are exactly x^k leaves, which correspond with the elements of the array. However, we will formulate our algorithm in such a way that it also works for trees with a non-uniform fan-out, for reasons explained in Section 4.

Every leaf l has an attribute $l.value : int$ that must be set to 1. The relation between the tree and the cube from Figure 3 is straightforward. Each level in the tree corresponds with a dimension, and a cell (c_0, c_1, c_2) of the cube corresponds with the leaf that we arrive at if we travel down the tree first taking the c_0 -th branch, then the c_1 -th branch, and finally the c_2 -th branch.

The internal nodes of the tree maintain information on how far the corresponding subtree has been processed already. Every internal node n has the following three attributes.

- $n.fan : int$
This constant denotes the number of children of the node.
- $n.nl : 0..n.fan$, initially 0
This variable denotes the number of child nodes that have already been processed, from left to right.
- $n.nr : 0..n.fan$, initially 0
This variable denotes the number of child nodes that have already been processed, from right to left.

Note that the subtree of node n has been processed completely if $n.nl + n.nr \geq n.fan$.

The root of the tree is denoted by *root* and the predicate *is_leaf* determines if a node is a leaf. Similar to algorithm *T* in [3], we make use of an atomic *compare-and-swap*-like instruction (see e.g. [5]). In the algorithm below this is denoted by placing angular brackets around the statement (‘*⟨*’ and ‘*⟩*’).

2.4 The algorithm

All processes operate in parallel and perform the same recursive procedure *traverse* with as the first argument the process identifier and the second argument the root of the tree. The recursive calls have as arguments smaller bit strings and other nodes of the tree. We use notation from [9] to express this.

```
forall pid in PID parbegin
  traverse(pid,root)
parend
```

Procedure *traverse* is defined below.

```
procedure traverse(bs,node)
var i: 0..node.fan;
begin
  if is_leaf(node) then
    node.value := 1
  else
    if head(bs) = 0 then
      i := node.nl;
      while i + node.nr < node.fan do
        traverse(tail(bs),child(node,i));
        ⟨ if node.nl = i then node.nl := i + 1 fi ⟩;
        i := node.nl
      od
    else
      i := node.nr;
      while node.nl + i < node.fan do
        traverse(tail(bs),child(node,node.fan - 1 - i));
        ⟨ if node.nr = i then node.nr := i + 1 fi ⟩;
        i := node.nr
      od
    fi
  fi
end
```

In the base case where the node is a leaf, the procedure writes the intended value in the array. Otherwise, the procedure treats the children of the node in a repetition from left to right or from right to left. The choice between starting left or right is irrelevant for correctness. For the sake of the complexity calculations, we let the choice depend on the head of the first argument *bs*, which is a suffix of the process’s *pid*. The recursive calls have the tail of the bit string *bs* as first argument, so that the processes start their actions at different points in the array. Private variable *i* is introduced to allow modification of the shared variables *node.nl* and *node.nr* by other processes.

It is worthwhile to notice that in the case that $N = P$ the above algorithm is equal to algorithm *X* (see [3]).

Since in this case we have $x = 2$ and the tree becomes a binary tree, which is traversed in exactly the same way as in algorithm *X*. The work calculations from Section 3.3 will show that in this case the upper bounds of the *Generalized Collision* algorithm and algorithm *X* are also identical.

3 Analysis of the algorithm

3.1 Correctness

The proof of correctness of the distributed algorithm consists of two steps. First, we prove partial correctness (i.e. if one of the processes successfully finishes, the whole tree has been processed) and, next, we prove termination (at least one process finishes successfully). If all leaves of a (sub)tree have been set to 1, we say that the (sub)tree has been processed.

Lemma 1. *The Generalized Collision algorithm is partially correct.*

Proof. Assuming that at least one of the processes finishes successfully, we have to prove that the whole tree has been processed. This follows immediately from the following two properties.

1. For every internal node *n* of the shared tree, it invariably holds that *n.nl* subtrees of node *n* from left to right have been processed. Likewise *n.nr* subtrees have been processed from right to left.
2. If a call *traverse(σ,n)* (for some bit string σ and some node *n*) finishes successfully, the subtree rooted in node *n* has been processed.

These two properties are proven with simultaneous induction on the depth of node *n*. The base case, where node *n* is a leaf is trivial. For the inductive case, we suppose that node *n* is an internal node.

1. The value of shared variable *n.nl* is only incremented from *i* to *i + 1* if procedure *traverse* has finished on the *i*th subtree of node *n*. By induction we then have that this subtree has been processed, which certifies this invariant. The variable *n.nr* is treated similarly.
2. If a call *traverse(σ,n)* finishes, one of the guards $i + n.nr < n.fan$ and $n.nl + i < n.fan$ must be false. Notice that in the first case we have $i \leq n.nl$ and in the second case $i \leq n.nr$. This is due to the fact that the values of *n.nr* and *n.nl* are non-decreasing. Therefore, if the call *traverse(σ,n)* finishes we have $n.nl + n.nr \geq n.fan$. Using the induction hypothesis, we can conclude that all subtrees of node *n* are processed, so the tree rooted in *n* is processed. \square

Next, we will prove termination of the algorithm.

Lemma 2. *The Generalized Collision algorithm terminates, i.e. at least one of the processes finishes successfully.*

Proof. For this we formulate the following termination function:

$$\sum_{n \in \text{internalnodes}} n.nl + n.nr$$

The fact that this is a proper termination function follows from the following observations.

First, the function is bounded. Since for every internal node n the values of $n.nl$ and $n.nr$ are bounded by the constant $n.fan$, the function is bounded by $2N^2$ (which is not a tight bound, see Section 3.3). Second, recall that the fault model implies that all but one process may fail. Since there are no blocking statements, the surviving process will continually invoke calls to procedure *traverse*, as long as it is not finished. After every call of this procedure (to, say, node n), the value of $n.nl + n.nr$ is strictly larger than before this call. Namely, it is incremented with 1 by the calling process, or it is incremented with at least 1 by one or more other processes. \square

In conclusion, we have that at least one call of *traverse(pid, root)* finishes successfully (termination) and that this implies that the complete tree rooted in *root* has been processed.

Corollary 1. *The Generalized Collision algorithm solves the asynchronous Write-All problem.*

3.2 Space usage

We will show that the processes have only moderate space requirements.

Lemma 3. *The space complexity of the Generalized Collision algorithm is $\mathcal{O}(N \log N + P \log^2 N)$.*

Proof. The shared data structure consists of the given array of N bits, together with the data at the internal nodes of the tree (see e.g. [9] for a description of how to represent a tree in a heap without overhead). There are less than N internal nodes. Every internal node n holds two shared variables of size $\log n.fan$. So the shared memory has size of $\mathcal{O}(N \log N)$.

Every process needs a private data structure with space for k stack frames, since the recursion depth is k . Each stack frame holds a local variable of size $\log n.fan$ and two parameters of sizes k and $\log N$. Since k is of order $\log N$, each process needs memory of order $\log^2 N$. \square

3.3 Work complexity

As was mentioned before, the work complexity of a parallel algorithm is the worst case total amount of work performed by the processes involved. With ‘total amount of work’ one generally means the number of instructions executed by all processes. We measure the work by counting the total number of calls of procedure *traverse* in a worst case scenario. The program text clearly shows that the number of instructions executed per call of *traverse*

is bounded by a constant. Therefore, the total number of procedure calls is an appropriate estimate here.

In the calculations below we will assume that the number of processes is 2^k and that the length of the array is x^k for some $k \geq 1$ and $x \geq 2$. This allows for the construction of a tree with a uniform fan-out. We will briefly consider the case of a tree with non-uniform fan-out in Section 4.

Because of the recursive structure of the input of the algorithm, the shared tree with fan-out x , we define the work inductively, i.e. express the work associated with a tree of height i in terms of the work associated with its subtrees of height $i - 1$. Note, that the number of processes, which of course plays an important role in determining the work, is 2^i (given a tree of height i). In our first inductive definition of work, however, we will decouple the number of processes and the height of the input tree, because, as we will see, subtrees of the input tree can be overloaded with processes (and *will* be overloaded in a worst case scenario).

We introduce $W_{i,j}$ as an upper bound on the work on a tree of height $i \leq j$ for a subsystem of 2^j processes with *pids* uniformly distributed for i . Here, we use the definition that a subsystem of 2^j processes has *pids* uniformly distributed for i if every bit sequence of length i is a suffix of the *pids* of precisely 2^{j-i} processes of the subsystem.

Lemma 4. *The work estimates $W_{i,j}$ satisfy the following recursive equations.*

$$W_{0,j} = 2^j \tag{1}$$

$$W_{i+1,j} = 2^j + (x - 1)W_{i,j-1} + W_{i,j} \tag{2}$$

Proof. Equation (1) is justified by the observation that, when 2^j processes start to work on a tree of height 0, a single leaf, they will all call procedure *traverse* once to set the array item associated with the leaf to 1, resulting in a work of 2^j .

Equation (2) is proved in the following way (see also Figure 4). When 2^j processes with *pids* uniformly distributed for $i + 1$ treat a tree of height $i + 1$, all of them first have to call procedure *traverse* at the root of the subtree. This accounts for the summand 2^j in the right-hand side of (2). Next, the processes split up in two groups of processes, according to the bit at position $i + 1$ from the end of their *pids*. Since the *pids* are uniformly distributed over $i + 1$, both groups have size 2^{j-1} and have *pids* uniformly distributed over i . One group treats the subtrees from left to right and the other group vice versa. The collision principle implies that the two groups interfere in at most one of the x subtrees, the one where they collide (the shaded sub-tree in Figure 4). The work associated with the subtrees can therefore be split in $x - 1$ times the work of 2^{j-1} processes on a tree of height i (the summand $(x - 1)W_{i,j-1}$), and the work of 2^j processes on a single tree of height i (the summand $W_{i,j}$). \square

In order to be able to simplify the recurrence relation from Lemma 4, we need the following property which states that doubling the number of processes doubles the work:

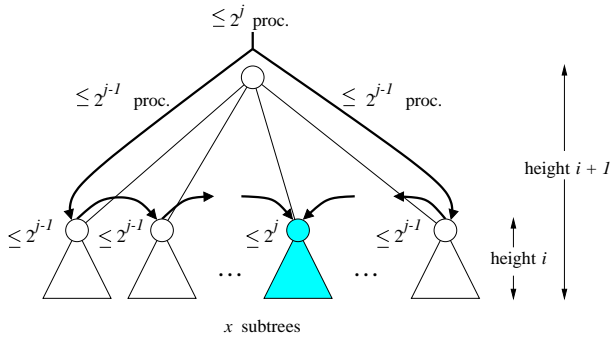


Fig. 4. Worst case distribution of processes over subtrees.

Lemma 5. For $j > i$, we have

$$2 \cdot W_{i,j-1} = W_{i,j} \quad (3)$$

Proof. We prove this by induction on i . For $i = 0$ we have $2 \cdot W_{0,j-1} = 2 \cdot 2^{j-1} = 2^j = W_{0,j}$. Assuming that the property holds for i , we derive for $i+1$: $2 \cdot W_{i+1,j-1} = 2 \cdot (2^{j-1} + (x-1) \cdot W_{i,j-2} + W_{i,j-1}) = 2^j + (x-1) \cdot 2 \cdot W_{i,j-2} + 2 \cdot W_{i,j-1} = 2^j + (x-1)W_{i,j-1} + W_{i,j} = W_{i+1,j}$. \square

The above property can also be explained in terms of *pids*. When the number of processes is doubled, they will have to share *pids* (because $j > i$). Each process will have a doppelgänger that follows the exact same route through the tree. This imitative behaviour explains the doubling of the work.

This leads to the introduction of w_i , which, for $i \geq 0$, denotes the work of 2^i processes on a tree of height i , so $w_k = W_{k,k}$.

Lemma 6. The work estimates w_k satisfy the following recursive equations.

$$w_0 = 1 \quad (4)$$

$$w_{i+1} = 2^{i+1} + (x+1)w_i \quad (5)$$

Proof. Because of Lemma 5 we can rewrite the second equation of the definition of W as follows:

$$\begin{aligned} W_{i+1,j} &= 2^j + (x-1)W_{i,j-1} + 2 \cdot W_{i,j-1} \\ &= 2^j + (x+1)W_{i,j-1} \end{aligned}$$

The equations now follow easily by setting $j = i + 1$. \square

Finally, we calculate an appropriate upper bound for w_k .

Theorem 1. The Generalized Collision algorithm solves the asynchronous *Write-All* problem with work complexity $\mathcal{O}(N \cdot P^{\log(\frac{x+1}{x})})$, where $x = N^{1/\log(P)}$.

Proof. Starting from the equations of Lemma 6, we arrive at the desired result with straightforward calcula-

tions.

$$\begin{aligned} w_k &= \{ \text{solve the recurrence relation} \} \\ &= \sum_{i=0}^k 2^{k-i} (x+1)^i \\ &= \{ \text{simple math} \} \\ &= \frac{(x+1)^{k+1} - 2^{k+1}}{x-1} \\ &= \{ \text{even more simple math} \} \\ &= \frac{x+1}{x-1} (x+1)^k - \frac{2^{k+1}}{x-1} \\ &\leq \{ x \geq 2, \text{ hence } \frac{x+1}{x-1} \leq 3 \text{ and } \frac{2^{k+1}}{x-1} > 0 \} \\ &= 3 \cdot (x+1)^k \end{aligned}$$

Hence, we have $w_k = \mathcal{O}((x+1)^k)$. We would like to express the work in terms of N and P , so we do some more calculations on $(x+1)^k$, using the equalities $N = x^k$ and $P = 2^k$:

$$\begin{aligned} (x+1)^k &= x^k \frac{(x+1)^k}{x^k} \\ &= x^k \left(\frac{x+1}{x} \right)^k \\ &= N \cdot \left(\frac{x+1}{x} \right)^{\log(P)} \\ &= N \cdot P^{\log(\frac{x+1}{x})} \end{aligned}$$

\square

4 Non-uniform fan-out

Although we were able to prove correctness of the algorithm for trees with non-uniform fan-out, we did assume uniform fan-out for our complexity calculations. This assumption proved very useful for obtaining a result which can easily be compared with work calculations for other algorithms.

Nonetheless, we claim that this assumption is not critical for the performance of our algorithm. Calculations and experimentation support this claim. Examples show that, strictly speaking, an optimal work complexity is almost never achieved with a uniform fan-out. In almost all cases the work complexity can be slightly improved by rebalancing the tree, while still keeping it quasi-uniform. By quasi-uniform we mean that the nodes at the same level of a tree have equal fan-out.

In the case of quasi-uniform fan-out the work load can be given as a closed expression that contains sum and product quantifiers. This goes as follows. Let us assume that every node at level i has fan-out x_i . The nodes at level 0 are leaves. So we have $x_0 = 0$. The analysis of Section 3.3 can be repeated and then yields instead of

formula (5) the recursive equation

$$w_{i+1} = 2^{i+1} + (x_{i+1} + 1)w_i$$

This implies that the total work load satisfies

$$w_k = \sum_{i=0}^k 2^i \prod_{j=i+1}^k (x_j + 1)$$

We now have to minimize the value of w_k under the constraint $\prod_{j=1}^k x_j \geq N$ for given value of k . It is not hard to make a functional program, e.g., in the language Haskell, to solve this optimization problem for given values of N and k .

It is even possible to find an approximate solution by analytic means. For that purpose, we define the real functions $f(x)$ and $g(x)$ with $x = (x_1, \dots, x_k) \in \mathbb{R}^k$ by

$$f(x) = \sum_{i=0}^k 2^i \prod_{j=i+1}^k (x_j + 1) \quad , \quad g(x) = \prod_{j=1}^k x_j - N$$

We are only interested in vectors x with all coordinates $x_i > 0$ (and preferably natural). So, now we have to minimize $f(x)$ under the constraint $g(x) = 0$ and all $x_i > 0$. According to the method of Lagrange multipliers (see e.g. [2] p. 315), we have that, if function f has an extremum at x under the constraint $g(x) = 0$, the gradient of f at x is a real multiple of the gradient of g . In this way, one finds that there is one optimal vector x and that it satisfies the recurrence relation

$$x_r = x_1 / (1 + 2 \sum_{i=1}^{r-1} \prod_{j=2}^i 2(x_j + 1)^{-1})$$

Note that this formula is independent of k . It implies that the numbers x_r form a decreasing sequence of positive reals. In particular, we have $x_2 = \frac{1}{3}x_1$. Of course, this only yields the optimum if the numbers x_i are allowed to be real.

The Haskell program mentioned above shows that for $N = 12000$ and $k = 5$, the optimal work load is obtained for the sequence x with $x_1 = 16$, $x_2 = 6$, and $x_3 = x_4 = x_5 = 5$. In the general case we see that an optimal work complexity is obtained if the fan-out for all levels are approximately the same, except for the fan-out at the level above the leaves, which should be three times larger.

We expect that in practice rebalancing the tree will yield at most a constant speed up in performance.

5 Observations and conclusions

We have presented an algorithm for the asynchronous *Write-All* problem. This algorithm is suitable for a multiprocess environment. It has good performance due to the lack of explicit synchronization. In particular this is the case when the task of setting a variable to one is replaced by a more time consuming operation. Moreover,

the algorithm is fault tolerant in the sense that it works correctly even if individual processes can fail or can stop and resume arbitrarily, assuming that not all processes die. Finally, our algorithm performs a kind of dynamic load balancing. Every process checks in a specific order all the tasks that must be executed and if it finds one that has not been performed, it carries it out. Due to the data structures involved, this can be done with minimal duplication of work. This guarantees a distribution of tasks over processes, where no process will idle when work can be done.

A potential drawback of our algorithm is that it utilizes on compare and swap registers. An interesting question is whether these can be replaced by atomic reads and writes. We believe that correctness of the algorithm is maintained by replacing the compare and swap register in a straightforward way by atomic reads and writes, but that the work increases. We believe that this is even the case when the compare and swap register is replaced by a test and set register.

Our algorithm improves upon existing asynchronous algorithms in several ways. In comparison with most published algorithms it has a better order of performance. This does not hold for algorithms AW and AW^T , which are based on a rather different algorithmic concept than our algorithm. Algorithm AW ‘only’ improves upon our algorithm with high probability, although we expect that in practice this algorithm has a good performance. From a theoretical perspective AW^T performs better than our algorithm, but due to a high initial constant amount of work AW^T is not suitable for any practical purposes.

To ascertain these findings, we have implemented our algorithm and have run it for different numbers of processes, where we compared the number of process steps with the worst case estimate of the amount of work that needs to be done. Without going into detail, as we believe that it is very hard to draw universal conclusions from experiments, we found that the work always remained far below our worst case estimate.

Finally, we make some observations concerning the restrictions on the values for N and P . In the case that we use a tree with uniform fan-out as the shared data structure, an array of size $N = x^k$ can be accommodated. However, such uniform fan-out is not needed for obtaining an optimal work complexity. By adjusting the fan-out of the nodes in the tree, it is possible to accommodate an array with arbitrary size N . Furthermore, since processes need not execute, we can take $P \leq 2^k$, provided all process identifiers differ and have a length at least equal to the depth of the tree. The work remains essentially the same.

Acknowledgements. We thank Dragan Bošnački, André Engels, Peter Hilbers, Jan Jongejan, Izak van Langevelde, Johan Lukkien, Alex Shvartsman, and the anonymous referees for comments, ideas and references.

References

1. R.J. Anderson and H. Woll. Algorithms for the certified write-all problem. *Siam Journal of Computing*, 26(5):1277-1283, 1997.
2. T.M. Apostol. *Calculus*, Vol. II. Wiley 1969.
3. J.F. Buss, P.C. Kanellakis, P.L. Ragde and A.A. Shvartsman. Parallel Algorithms with Processor Failures and Delays. *Journal of Algorithms*, 20:45-86, 1996.
4. J.F. Groote and W.H. Hesselink. Synchronization-free parallel accessible hash-tables. In preparation, 2000.
5. M.P. Herlihy. Wait-free synchronization. *ACM Trans. on Program. Languages and Systems* **13** (1991) 124–149.
6. W.H. Hesselink and J.F. Groote. Waitfree Distributed Memory Management by Create, and Read Until Deletion (CRUD). Technical report SEN-R9811, CWI, Amsterdam, 1998.
7. P.C. Kanellakis and A.A. Shvartsman. Efficient parallel algorithms can be made robust. *Distributed Computing*, 5(4):201–217, 1992. A preliminary version appeared in *Proceedings of the 8th ACM PODC*, pages 211–222, 1989.
8. P.C. Kanellakis and A.A. Shvartsman. Fault-tolerance and efficiency in massively parallel algorithms. In G.M. Koob and C.G. Lau, editors, *Foundations of Dependable Computing – Paradigms for Dependable Applications*, pages 125–154, Kluwer Academic, 1994.
9. P.C. Kanellakis and A.A. Shvartsman. *Fault-tolerant parallel computation*. Kluwer Academic Publishers, 1997.
10. Z.M. Kedem, K.V. Palem, A. Raghunathan, and P. Spirakis. Combining tentative and definite executions for dependable parallel computing in *Proceedings of the 23rd ACM Symposium on Theory of Computing*, 1991.