

# 1 Regular Tree Algorithm Redux: The Atomic Case: 2 DREU Final Report

3 **Soma Chaudhuri**  
4 Iowa State University, USA  
5 chaudur@iastate.edu

6 **Reginald Frank**  
7 Texas A&M University, USA  
8 reginaldfrank77@tamu.edu

9 **Jennifer L. Welch**<sup>1</sup>  
10 Texas A&M University, USA  
11 welch@cse.tamu.edu

## 12 — Abstract —

13 We consider the problem of simulating a  $k$ -valued atomic register in a wait-free manner using  
14 binary atomic registers as building blocks, for any  $k > 2$ . In this note we prove that the tree  
15 algorithm of Chaudhuri and Welch [1], which was originally proposed to simulate a  $k$ -valued  
16 *regular* register out of binary *regular* registers, simulates an *atomic* register when the binary  
17 registers are *atomic*. Our tree algorithm uses  $k - 1$  binary registers and each operation on the  
18 simulated register performs  $\lceil \log_2 k \rceil$  operations on the binary registers. When the number of  
19 readers of the simulated register is  $\omega\left(\frac{k}{\log k}\right)$ , this is the most space-efficient algorithm known for  
20 the problem. Furthermore we prove that the said algorithm's worst case logical WRITE use the  
21 fewest amount of physical writes possible.

22 **2012 ACM Subject Classification** Theory of computation → Distributed algorithms

23 **Keywords and phrases** Interprocess communication, Shared registers

## 24 **1** Introduction

25 A *register* is a fundamental shared object which can either be read from or written to. Our  
26 work focuses on registers with a single writer, but an arbitrary number of readers. Different  
27 registers may have different properties such as *regularity* or *atomicity*. *Regularity* ensures  
28 that if a read of the register is concurrent with a write, then the read will return either the  
29 previous (if no preceding writes, initial) value of the register or the value of the concurrent  
30 write; otherwise the read will return the most recent value held/written to the register.  
31 *Atomicity* is stronger than regularity as it provides a definite order to the operations in an  
32 execution, meaning the operations are *linearizable* [4]. Furthermore it is accepted [3] that  
33 the only distinction between regular and atomic executions is that regular executions allow  
34 for *new-old* inversions. These inversions occur when one or more read returns an *old* value,  
35 the value which was written just before the current write after the current write has been  
36 read. Reads of a regular register are permitted to return old values allowing for *new-old*  
37 inversions, however atomicity forbids this.

38 Several algorithms have been created to make use of a number of binary registers for the  
39 purpose of constructing  $k$ -valued registers, a register which can hold  $k$  unique values where  
40  $k \in \mathbb{Z}^+$ . In particular, the tree algorithm created by Chaudhuri and Welch in [1] uses  $k - 1$

---

<sup>1</sup> Supported in part by NSF grant 1526725.

41 binary regular registers to construct a  $k$ -valued regular register. The motive of this document  
 42 is to show that this algorithm can directly be extended to use atomic binary registers to  
 43 create an atomic  $k$ -valued register. This will be done by showing that “new-old inversions”  
 44 are not possible in the resulting construction.

## 45 **2** Tree Algorithm

46 The original tree algorithm presented in [1] goes as follows. For a value set  $V$  of size  $k$ ,  
 47 associate a binary register with each internal node of a binary tree with  $k$  leaves and label  
 48 the leaves with the elements of  $V$ . A READ operation reads the registers corresponding to a  
 49 path from the root down to a leaf. At each level, if the current register holds 0 (resp., 1),  
 50 then the next register read is the left (resp., right) child. If the last register in the path  
 51 holds 0 (resp., 1), then return the value labeling the left (resp., right) child of the tree node  
 52 corresponding to the register. A WRITE( $v$ ) operation writes the registers corresponding to  
 53 the path from the leaf labeled  $v$  up to the root. At each level, if the previous tree node is  
 54 the left (resp., right) child of the current tree node, then 0 (resp., 1) is written to the register  
 55 corresponding to the current tree node. The tree algorithm uses  $k - 1$  binary registers. If  
 56 the binary tree is complete, then the number of steps in each operation is at most  $\lceil \log_2 k \rceil$ .

57 It is shown in [1] that if the binary registers are regular, then the  $k$ -valued register  
 58 simulated by the tree algorithm is also regular. We will show that if the binary registers are  
 59 atomic, then the tree algorithm simulates an atomic register.

60 To simplify the correctness proof, we provide a new way of describing the tree algorithm,  
 61 called  $RTA(V)$  for Recursive Tree Algorithm, where  $V$  is the finite set of ordered values to be  
 62 stored in the register. For now, assume that  $|V| = k$  is a power of 2. (Later we will discuss  
 63 how to relax this assumption.)

64 For the base case, when  $|V| = 2$ ,  $RTA(V)$  uses a single binary<sup>2</sup> atomic register  $x$ , the  
 65 READ algorithm consists solely of reading  $x$  and returning the value obtained, and the  
 66 WRITE( $v$ ) algorithm consists solely of writing  $v$  to  $x$ .

67 Now suppose  $|V| > 2$ .  $RTA(V)$  can be viewed as a binary tree with two levels. The root  
 68 of the tree is a binary atomic register, *root*, with two children. The left child, denoted as **A**,  
 69 and the right child, denoted as **B**, are both  $\frac{k}{2}$ -valued atomic registers, where **A**'s value set  
 70 consists of the  $\frac{k}{2}$  smallest values in  $V$  and **B**'s value set consists of the  $\frac{k}{2}$  largest values in  $V$ .

71 The READ algorithm first reads the root; if 0 (resp., 1) is returned, then it reads **A** (resp.,  
 72 **B**) and returns the value obtained. Pseudocode is in Algorithm 1. To WRITE  $v$  when  $v$  is  
 73 less than the median of  $V$ , the algorithm first writes  $v$  to **A** and then writes 0 to the root;  
 74 when  $v$  is greater than the median of  $V$ , it first writes  $v$  to **B** and then writes 1 to the root.  
 75 Pseudocode is in Algorithm 2.

---

<sup>2</sup> The size of the value set is 2, but the values are not necessarily 0 and 1.

76

**Algorithm 1** *RTA READ*,  $|V| > 2$ 


---

```

1: procedure READ()
2:    $rootVal \leftarrow \text{read}(root)$ 
3:   if  $rootVal = 0$  then
4:     RETURN read(A)
5:   else
6:     RETURN read(B)

```

---

77

78 To simulate the  $k$ -valued register from binary registers, use *RTA* recursively in this  
79 construction. That is, let  $\mathbf{A} = RTA(V_1)$ , where  $V_1$  consists of the  $\frac{k}{2}$  smallest elements of  $V$   
80 and let  $\mathbf{B} = RTA(V_2)$ , where  $V_2$  consists of the  $\frac{k}{2}$  largest elements of  $V$ .

81 When the recursion is unrolled, the resulting algorithm is equivalent to the tree algorithm  
82 in [1]. Informally the reason is that all key features of the tree algorithm, such as WRITES  
83 progressing from leaf to root, and READS progressing from root to leaf, are preserved.

**3 Analysis**

84  
85 In this section, we show that  $RTA(V)$  simulates an atomic register.

86 For an execution of an atomic register  $x$ , we say that read  $r$  of  $x$  reads from a write  $w$  of  
87  $x$  if  $w$  is the latest write that precedes  $r$  in the linearization of the operations on  $x$ .

88 **► Theorem 1.** *RTA(V) implements an atomic register when  $|V| = k$  is a power of 2.*

89 **Proof.** We use induction on the size of  $V$  to show that  $RTA(V)$  is atomic.

90 *Base Case:*  $|V| = 2$ .  $RTA(V)$  is simply a binary atomic register and the result follows.

91 *Inductive Step:*  $|V| = 2^m$ , where  $m > 1$ . By the inductive hypothesis, the building block  
92 registers  $\mathbf{A}$  and  $\mathbf{B}$  are atomic since each of them is an instantiation of *RTA* for a value set of  
93 size  $2^{m-1}$ .

94 By a result of Lamport [4], to show that an execution is atomic, it is sufficient to  
95 demonstrate the existence of a function  $\rho$  from the set of READS in the execution to the set  
96 of WRITES in the execution such that:

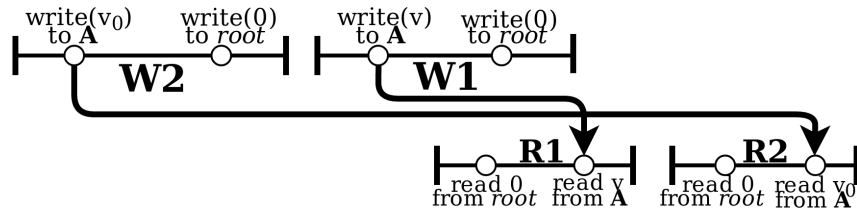
- 97 1. For every READ  $R$ ,  $\rho(R)$  either overlaps  $R$  or is the latest WRITE that finishes before  $R$   
98 begins.
- 99 2. For all READS  $R_1$  and  $R_2$  such that  $R_1$  finishes before  $R_2$  begins,  $\rho(R_1)$  either equals or  
100 precedes  $\rho(R_2)$ .

101 The first condition ensures regularity and the second condition rules out “new-old” inversions.

102 For the tree algorithm, we define the function  $\rho$  as follows. For READ  $R$ , let  $\rho(R)$  be the  
103 WRITE  $W$  such that  $R$ 's leaf read (of  $\mathbf{A}$  or  $\mathbf{B}$ ) reads from  $W$ 's leaf write. Note that the  
104 value written to the leaf is the same as the value read from that leaf, and thus the value  
105 returned by  $R$  is the same as the value written by  $W$ .

106 It is straightforward to see that  $\rho$  satisfies condition (1). Suppose for contradiction that  
107  $\rho$  does not satisfy condition (2). Then there are two READS,  $R_1(v_1)$  and  $R_2(v_2)$ , such that  
108  $R_1$  finishes before  $R_2$  begins but  $\rho(R_2) = W_2(v_2)$  finishes before  $\rho(R_1) = W_1(v_1)$  begins.

109 **Case 1:**  $R_1$  and  $R_2$  read the same value from  $root$ , without loss of generality, say 0. This  
110 means they both read leaf  $\mathbf{A}$ , resulting in the following precedence relations (denoted  $\rightarrow$ );  
111 see Figure 1.



■ **Figure 1** Precedence relations for Case 1

- 112 ■ read(0) from *root* in  $R_1 \rightarrow$  read( $v_1$ ) from **A** in  $R_1 \rightarrow$  read(0) from *root* in  $R_2 \rightarrow$  read( $v_2$ )
- 113 from **A** in  $R_2$
- 114 ■ write( $v_2$ ) to **A** in  $W_2 \rightarrow$  write(0) to *root* in  $W_2 \rightarrow$  write( $v_1$ ) to **A** in  $W_1 \rightarrow$  write(0) to
- 115 *root* in  $W_1$
- 116 ■ write( $v_1$ ) to **A** in  $W_1 \rightarrow$  read( $v_1$ ) from **A** in  $R_1$
- 117 ■ write( $v_2$ ) to **A** in  $W_2 \rightarrow$  read( $v_2$ ) from **A** in  $R_2$

118 Since  $W_2$ 's write to **A** is followed by  $W_1$ 's write to **A**, which is followed by  $R_1$ 's read of **A**,  
 119 which is followed by  $R_2$ 's read of **A**, it is not possible for  $\rho(R_2)$  to be  $W_2$ , a contradiction.

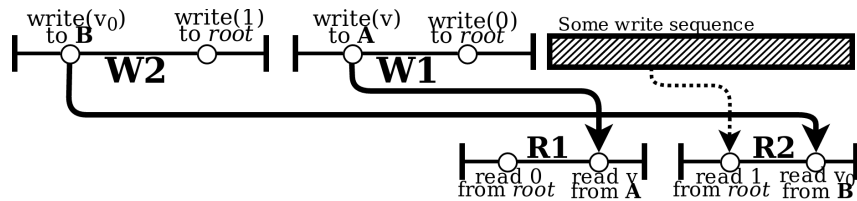
120 **Case 2:**  $R_1$  and  $R_2$  read different values from *root*; without loss of generality suppose  $R_1$   
 121 reads 0 and thus reads leaf **A**, while  $R_2$  reads 1 and thus reads leaf **B**.

122 We now have the following precedence relations:

- 123 ■ read(0) from *root* in  $R_1 \rightarrow$  read( $v_1$ ) from **A** in  $R_1 \rightarrow$  read(1) from *root* in  $R_2 \rightarrow$  read( $v_2$ )
- 124 from **B** in  $R_2$
- 125 ■ write( $v_2$ ) to **B** in  $W_2 \rightarrow$  write(1) to *root* in  $W_2 \rightarrow$  write( $v_1$ ) to **A** in  $W_1 \rightarrow$  write(0) to
- 126 *root* in  $W_1$
- 127 ■ write( $v_1$ ) to **A** in  $W_1 \rightarrow$  read( $v_1$ ) from **A** in  $R_1$
- 128 ■ write( $v_2$ ) to **B** in  $W_2 \rightarrow$  read( $v_2$ ) from **B** in  $R_2$

129 One important observation is that if a WRITE writes 0 (resp., 1) to *root* then it has  
 130 previously written to leaf **A** (resp., **B**).

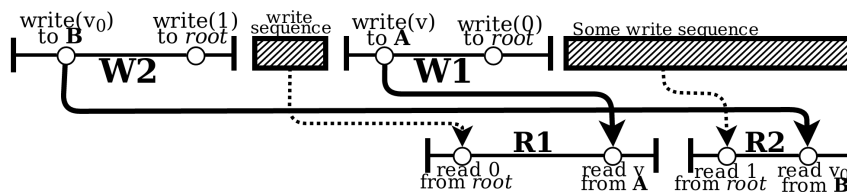
131 **Case 2.1:**  $R_1$ 's read of *root* reads from  $W_1$ 's write to *root* or a later one. See Figure 2.  
 132 Then  $R_2$ 's read of *root* reads from a write to *root* in  $W_1$  or a later WRITE. By the observation,  
 133  $R_2$ 's read of **B** reads from the write to **B** in  $W_1$  or a later WRITE, so  $\rho(R_2)$  cannot be  $W_2$ , a  
 134 contradiction.



■ **Figure 2** Precedence relations for Case 2.1

135 **Case 2.2:**  $R_1$ 's read of *root* reads from the write to *root* in a WRITE earlier than  $W_1$ . It  
 136 cannot be from  $W_2$ , since  $W_2$  writes a different value to *root* than what  $R_1$  reads. It cannot  
 137 be from a WRITE that precedes  $W_2$ , since earlier writes to *root* are followed by  $W_2$ 's. So it  
 138 must be from some WRITE  $W_3$  that is in between  $W_2$  and  $W_1$ . See Figure 3. But then  $W_2$ 's  
 139 write to *root* precedes  $W_3$ 's write to *root*, which precedes  $R_1$ 's read of *root*, which precedes  
 140  $R_2$ 's read of *root*. Thus  $R_2$ 's read of *root* reads from a write to *root* in some WRITE  $W_4$  that

141 follows  $W_3$ . By the observation,  $W_4$  writes  $\mathbf{B}$  before it writes  $root$ , and thus  $\pi(R_2)$  cannot  
 142 be  $W_2$ , a contradiction.



143 **Figure 3** Precedence relations for Case 2.2

143 Thus  $\rho$  satisfies condition 2 and hence  $RTA(V)$  is atomic.  $\blacktriangleleft$

144 As noted in [1] for the original tree algorithm,  $RTA(V)$  uses  $k - 1$  binary registers, where  
 145  $|V| = k$ , and the number of steps per operation is  $\log_2 k$ .

## 146 4 Extensions

147 Because  $RTA$  is equivalent to the Tree Algorithm with atomic binary registers the same  
 148 method can be used to allow for a value of  $k$  which is not a power of 2. The general  
 149 intuition for this is that each implementation where  $k$  is not a power of two,  $RTA$  provides  
 150 a tree which is equivalent to a tree where  $k$  is the following power of 2 with a number  
 151 of nodes removed. The restricted nodes has an equivalent purpose as restricting WRITES  
 152 to only use values between 0 and  $k$ ; it follows that the all the possible executions of this  
 153 restricted implementation are a subset of the set of possible executions when WRITES are  
 154 not restricted. Because the unrestricted set of executions are all linearizable, it follows then  
 155 that the restricted set is also linearizable showing that the original restriction on  $k$  is not  
 156 necessary to implement an atomic  $RTA$  algorithm.

157 This shows that  $RTA$  can correctly implement an atomic register for any sized  $k$ , allowing  
 158  $\mathbf{A}$  and  $\mathbf{B}$  to be created using  $RTA$  meaning all implementations of  $RTA$  can be constructed  
 159 solely using binary atomic registers.

## 160 5 Lower Bound

161 In this section we give several facts about algorithms where all logical READS use  $\leq \log_2(k)$   
 162 physical reads.

163 **Theorem 2.** For any wait-free algorithm  $A$  that simulates a  $k$ -valued safe register using  
 164 binary atomic registers, if the READ algorithm uses at most  $\log_2(k)$  reads in the worst case,  
 165 then the WRITE algorithm uses at least  $\log_2(k)$  writes in the worst case.

166 **Proof.** Consider the following set  $S$  of executions of  $A$ :  $\sigma_i = \text{WRITE}(v_i) \alpha_i \text{ ACK READ}_1 \beta_i$   
 167  $\text{RETURN}(w_i) \mid$  only the  $WP$  takes steps in  $\alpha_i$  and only  $RP_1$  takes steps in  $\beta_i$ ,  $0 \leq i < k$ .  
 168 Note that  $w_i$  must equal  $v_i$  to satisfy safety.

169 Now construct a decision tree of  $RP_1$ 's behavior in the executions in  $S$ . The root of the  
 170 decision tree corresponds to the first register that  $RP_1$  reads, this register is the same for all  
 171 executions in  $S$  as  $RP_1$  always begins in the same state. Now depending on the first value  
 172 read,  $RP_1$  does some amount of non-read actions then does another read of some register.  
 173 Because the root only can store two values there are only two registers  $RP_1$  can decide to  
 174 read from causing the root to have two children. Continue to build the tree like this. Add  
 175 leaves to indicate which value  $RP_1$  will RETURN.

176 ▶ **Lemma 3.** *The decision tree is a complete binary tree with exactly  $k$  leaves in which every*  
 177 *leaf is at depth  $\log_2(k)$  and the path from the root to the leaf labeled  $v_i$  corresponds to  $\beta_i \in \sigma_i$ .*  
 178

179 **Proof.** Note that there must be at least  $k$  leaves in the decision tree, as each value in  $V$  is  
 180 returned in one of the executions in  $S$ . Since the decision tree is binary, basic facts from  
 181 graph theory imply that each leaf must be at depth  $\log_2(k)$ , allowing exactly  $k$  leaves and  
 182 each root-to-leaf path must correspond to a different execution in  $S$ . ◀

183 ▶ **Lemma 4.** *None of the READs in an execution in  $S$  reads the same register more than*  
 184 *once.*

185 **Proof.** Suppose in contradiction the READ in  $\sigma_i$  reads some number of registers more than  
 186 once, for some  $i$ . Let  $x$  be the first register that  $RP_1$  reads twice, say the  $a$ -th read and the  
 187  $b$ -th read, with  $\log_2(k) \geq b > a$ . Consider the node in the decision tree for the  $b$ -th read  
 188 in the path for  $\sigma_i$ . This node must have two children since the decision tree is complete.  
 189 Consider a root-to-leaf path  $\pi$  in the decision tree that forks off from the path corresponding  
 190 to  $\sigma_i$  at this node. By Lemma 1, there is a one-to-one correspondence between root-to-leaf  
 191 paths in the decision tree and the set of executions and so  $\pi$  corresponds to  $\sigma_j$  in  $S$ , for some  
 192  $j \neq i$ .

193 However, it is not possible for  $\sigma_i$  and  $\sigma_j$  to read different values from  $x$  at the  $b$ -th read  
 194 since the only process that is taking steps during the READs is  $RP_1$ : even if  $RP_1$  writes to  $x$   
 195 during the READ, it will write the same value in  $\sigma_i$  as in  $\sigma_j$ , as nothing differs in those two  
 196 executions until reaching the  $b$ -th read. ◀

197 Now we can finish the proof of the theorem. Let  $x_1$  be the first register read in the READ  
 198 of each execution in  $S$ . Note that in half of the elements of  $S$ , the value read from  $x_1$  in  
 199 the READ must be different from the initial value of  $x_1$ , in order to be able to reach all the  
 200 leaves on that half of the decision tree. In other words, the WRITE writes to  $x_1$  in half  
 201 the executions in  $S$ . Let  $S_1$  be the subset of  $S$  consisting of the executions in which the  
 202 WRITE writes  $x_1$ ; note that  $|S_1| = \frac{k}{2}$ . Now let  $x_2$  be the register corresponding to the root  
 203 of subtree in which  $x_1$  was written. A similar argument shows that in half the executions  
 204 in  $S_1$ , the value read from  $x_2$  in the READ must be different from the initial value of  $x_2$ .  
 205 Let  $S_2$  be the subset of  $S_1$  consisting of the executions in which the WRITE writes to  $x_2$ ;  
 206 note that  $|S_2| = \frac{k}{2^2}$ . Continuing this way we obtain  $S_{\log_2 k}$ , which is of size 1, in which the  
 207 WRITE writes to  $x_1, x_2, \dots, x_{\log_2 k}$ . By Lemma 2, each of these registers is distinct and thus  
 208 the WRITE writes to at least  $\log_2(k)$  registers. ◀

## 209 6 Conclusion

210 This result directly improves on the Tree-Algorithm found in [2], however it does not improve  
 211 their combined implementation which uses the un-optimal tree algorithm. And in conclusion  
 212 the tree algorithm uses an optimal number steps in both READs and WRITES when all READs  
 213 use  $\leq \log_2(k)$  physical reads.

## 214 — References —

- 215 1 Soma Chaudhuri and Jennifer L. Welch. Bounds on the costs of multivalued register im-  
 216 plementations. *SIAM J. Comput.*, 23(2):335–354, 1994.

- 217 **2** Tian Ze Chen and Yuanhao Wei. Step optimal implementations of large single-writer  
218 registers. In *20th International Conference on Principles of Distributed Systems (OPODIS)*,  
219 pages 32:1–32:16, 2016.
- 220 **3** Rachid Guerraoui and Petr Kuznetsov. *Concurrent Computing*, January 21, 2018 (accessed  
221 July 4, 2018). <https://perso.telecom-paristech.fr/kuznetso/MITRO207-2018/book-ln.pdf>.
- 222 **4** Leslie Lamport. On interprocess communication. Part II: Algorithms. *Distributed Comput-*  
223 *ing*, 1(2):86–101, 1986.