

# Brief Announcement: Leader Election for Arbitrarily Connected Networks in the Presence of Process Crashes and Weak Channel Reliability

Sergio Rajsbaum<sup>†</sup>, Michel Raynal<sup>\*;‡</sup>, Karla Vargas<sup>†</sup>

<sup>†</sup>Instituto de Matematicas, UNAM, Mexico

<sup>\*</sup>Univ Rennes IRISA, 35042 Rennes, France

<sup>‡</sup>Department of Computing, Polytechnic University, Hong Kong

## 1 Introduction

In the leader election problem each process  $p_i$  has a local variable  $leader_i$ , and it is required that all the local variables  $leader_i$  forever contain the same identity, which is the identity of one of the processes. If processes may crash, the system is fully asynchronous, and the elected leader must be a process that does not crash, leader election cannot be solved [4]. Not only the system must no longer be fully asynchronous, but the leader election problem must be weakened to the *eventual leader election problem*. This problem is denoted  $\Omega$  in the failure detector parlance [1,2]. Notice that the algorithm must elect a new leader each time the previously elected leader crashes.

ADD channels were introduced in [5], as a realistic partially synchronous model of channels that can lose and reorder messages. Each channel guarantees that some subset of the messages sent on it will be delivered in a timely manner and such messages are not too sparsely distributed in time. More precisely, for each channel there exist two constants  $K$  and  $D$ , not known to the processes (and not necessarily the same for all channels), such that for every  $K$  consecutive messages sent in one direction, at least one is delivered within  $D$  time units after it has been sent.

Even though ADD channels seem so weak, it is possible to implement an *eventually perfect failure detector* in an arbitrarily connected network of ADD channels [3]. A implementation of  $\diamond P$  using messages of size  $O(n \log n)$  in the same model was presented [6]. The goal of this paper is move from  $\diamond P$  to  $\Omega$  using messages of  $O(\log n)$ .

This paper shows that it is possible to implement  $\Omega$  in an arbitrarily connected network of eventual ADD channels where asynchronous processes may fail by crashing using messages of  $O(\log n)$ . Then, we propose an implementation of  $\Omega$  in networks with unknown membership whose messages are eventually of size  $O(\log n)$  too.

Designing leader election ADD-based algorithms using messages whose size is bounded, is a difficult challenge since while the constants  $K$  and  $D$  do exist. We found it even more challenging to work under the assumption that some edges might not satisfy any property at all; our algorithm works under the assumption that only edges on an (unknown to the processes) spanning tree are guaranteed to comply with the ADD property.

## 2 Model of Computation

The system consists of a finite set of processes  $\Pi = \{p_1, p_2, \dots, p_n\}$ . Any number of processes may fail by crashing. A process is *correct* if it does not crash, otherwise, it

is *faulty*. The communication network is represented by a directed graph  $G = (II, E)$ , where an edge  $(p_i, p_j) \in E$  means that there is a unidirectional channel that allows the process  $p_i$  to send messages to  $p_j$ . It is required the existence of a spanning tree containing all correct processes and the root being the leader, i.e. the correct process with the smallest identity.

A directed channel  $(p_i, p_j)$  satisfies the *ADD property* if there are two constants  $K$  and  $D$  (unknown to the processes such that for every  $K$  consecutive messages sent by  $p_i$  to  $p_j$ , at least one is delivered to  $p_j$  within  $D$  time units after it has been sent. The other messages from  $p_i$  to  $p_j$  can be lost or experience arbitrary delays.

|   |                   |
|---|-------------------|
| <b>initialization</b>   | —Code for $p_i$ — |
| (1) $leader_i \leftarrow i$ ; $hopbound_i[i] \leftarrow n$ ; set $timer_i[i, n]$ to $+\infty$ ;   |                   |
| (2) <b>for each</b> $j \in \{1, \dots, n\} \setminus \{i\}$ <b>and each</b> $x \in \{1, \dots, n\}$   |                   |
| (3) <b>do</b> $timeout_i[j, x] \leftarrow$ any positive integer; set $timer_i[j, x]$ to $timeout_i[j, hb]$ ;                                  |                   |
| (4)     set $penalty_i[j, x]$ to $-1$ ; $hopbound_i[j] \leftarrow 0$  |                   |
| (5) <b>end for.</b>   |                   |
| (6) <b>every</b> $T$ time units of $clock_i()$ <b>do</b>  |                   |
| (7) <b>if</b> ( $hopbound_i[leader_i] > 1$ )  |                   |
| (8) <b>then for each</b> $j \in out\_neighbors_i$   |                   |
| <b>do</b> send ALIVE( $leader_i, hopbound_i[leader_i] - 1$ ) to $p_j$ <b>end for</b>  |                   |
| (9) <b>end if.</b>  |                   |
| (10) <b>when</b> ALIVE( $\ell, hb \leftarrow n - k$ ) <b>such that</b> $\ell \neq i$ <b>is received</b> % from a process in $in\_neighbors_i$ |                   |
| (11) <b>if</b> ( $\ell \leq leader_i$ )   |                   |
| (12) <b>then</b> $leader_i \leftarrow \ell$ ;   |                   |
| (13) <b>if</b> ( $[timer_i[leader_i, hb]$ expired)  |                   |
| <b>then</b> increase the value of $timeout_i[leader_i, hb]$ <b>end if</b> ;   |                   |
| (14)         set $timer_i[leader_i, hb]$ to $timeout_i[leader_i, hb]$ ;   |                   |
| (15) $not\_expired_i \leftarrow \{x \mid timer_i[leader_i, x] \text{ not expired } \}$ ;  |                   |
| (16) $hopbound_i[leader_i] \leftarrow$  |                   |
| $\max\{x \in not\_expired \text{ with smallest non-negative } penalty_i[leader_i, x]\}$   |                   |
| (17) <b>end if.</b>   |                   |
| (18) <b>when</b> $timer_i[leader_i, hb]$ <b>expires</b> and ( $leader_i \neq i$ ) <b>do</b>   |                   |
| (19) $penalty_i[leader_i, hb] \leftarrow penalty_i[leader_i, hb] + 1$ ;   |                   |
| (20) <b>if</b> ( $\wedge_{1 \leq x \leq n} ([timer_i[leader_i, x]$ expired))  |                   |
| (21) <b>then</b> $leader_i \leftarrow i$  |                   |
| (22) <b>else</b> same as lines 15-16  |                   |
| (23) <b>end if.</b>   |                   |

Algorithm 1: Eventual leader election in the  $\diamond$ ADD model with known membership

### 3 An Algorithm for Eventual Leader Election in the $\diamond$ ADD Model with Unknown Membership

The second algorithm (only described here due space limitations) solves eventual leader election in the  $\diamond$ ADD model with unknown membership, which means that, initially, a process knows nothing about the network, it knows only its input/output channels.

Initially  $p_i$  communicate its identity to its neighbors. Once its neighbors know about it,  $p_i$  no longer send its identity. And the same with other names that  $p_i$  learns. For that,  $p_i$  keeps a *pending set* for every channel connected to it that helps it to keep track of the information that it needs to send to its neighbors. So initially,  $p_i$  adds the pair  $(\text{new}, i)$  to every pending set.

When process  $p_i$  receives an `ALIVE()` message from  $p_j$ , this message can contain information about the leader and the corresponding pending set that  $p_j$  saves for  $p_i$ . First,  $p_i$  processes the information contained in the pending set and then processes the information about the leader. If  $p_i$  finds a pair with a name labeled as `new` and does not know it, it stores the new name in the set  $knnonw_i$ , increases its hopbound value, and adds to every pending set (except to the one belonging to  $p_j$ ) this information labeled as `new`. In any case,  $p_i$  needs to communicate  $p_j$  that it already knows that information, so  $p_i$  adds this information to the pending set of  $p_j$  but labeled as an acknowledgment.

When  $p_j$  receives `name` labeled as an acknowledgment from  $p_i$ , i.e.  $(\text{ack}, \text{name})$ , it stops sending the pair  $(\text{new}, \text{name})$  to it, so it deletes that pair from  $p_i$ 's pending set. Eventually,  $p_i$  receives a pending set from  $p_j$  not including  $(\text{new}, \text{name})$ , so  $p_i$  deletes  $(\text{ack}, \text{name})$  from  $p_j$ 's pending set.

As in Algorithm 1, every process keeps as leader a process with minimum id. This part is similar to Algorithm 1, only ignoring the penalties since we are assuming that all channels are  $\diamond$  ADD.

### 4 Conclusion

The  $\diamond$ ADD model is a particularly weak partially synchronous communication model. Assuming first that the correct processes are connected by a spanning tree made up of  $\diamond$  ADD channels, this article has presented an algorithm that elects an eventual leader, using messages of only size  $O(\log n)$ .

### References

1. Chandra T.D., Hadzilacos V., and Toueg S., The weakest failure detector for solving consensus. *Journal of the ACM*, 43(4):685-722 (1996)
2. Chandra T.D. and Toueg S., Unreliable failure detectors for reliable distributed systems. *Journal of the ACM*, 43(2):225-267 (1996)
3. Kumar S. and Welch J.L., Implementing  $\diamond P$  with bounded messages on a network of ADD channels. *Parallel Processing Letters*, Vol. 29(1) 1950002, 12 pages (2019)
4. Raynal M., *Fault-tolerant message-passing distributed systems: an algorithmic approach*. Springer, 492 pages, ISBN 978-3-319-94140-0 (2018)
5. Sastry S. and Pike S.M., Eventually perfect failure detectors using ADD channels. *5th Int'l Symposium on Parallel and Distributed Processing and Applications (ISPA'07)*, Springer LNCS 4742, pp. 483-496 (2007)
6. Vargas K. and Rajsbaum S., An eventually perfect failure detector for networks of arbitrary topology connected with ADD channels using time-to-live values. *49th IEEE/IFIP Int'l Conference on Dependable Systems and Networks (DSN'19)*, IEEE Press, pp. 264-275 (2019)