

# Data redistribution algorithms for heterogeneous processor rings

Hélène Renard, Yves Robert, Frédéric Vivien

► **To cite this version:**

Hélène Renard, Yves Robert, Frédéric Vivien. Data redistribution algorithms for heterogeneous processor rings. [Research Report] Laboratoire de l'informatique du parallélisme. 2004, 2+26p. hal-02102059

**HAL Id: hal-02102059**

**<https://hal-lara.archives-ouvertes.fr/hal-02102059>**

Submitted on 17 Apr 2019

**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.



**Laboratoire de l'Informatique du Parallélisme**

École Normale Supérieure de Lyon  
Unité Mixte de Recherche CNRS-INRIA-ENS LYON-UCBL n° 5668

***Data redistribution algorithms  
for heterogeneous processor rings***

Hélène Renard,  
Yves Robert,  
Frédéric Vivien

May 2004

Research Report N° 2004-28

**École Normale Supérieure de Lyon**

46 Allée d'Italie, 69364 Lyon Cedex 07, France

Téléphone : +33(0)4.72.72.80.37

Télécopieur : +33(0)4.72.72.80.80

Adresse électronique : [lip@ens-lyon.fr](mailto:lip@ens-lyon.fr)



# Data redistribution algorithms for heterogeneous processor rings

Hélène Renard, Yves Robert, Frédéric Vivien

May 2004

## Abstract

We consider the problem of redistributing data on homogeneous and heterogeneous ring of processors. The problem arises in several applications, each time after that a load-balancing mechanism is invoked (but we do not discuss the load-balancing mechanism itself). We provide algorithms that aim at optimizing the data redistribution, both for uni-directional and bi-directional rings, and we give complete proofs of correctness. One major contribution of the paper is that we are able to prove the optimality of the proposed algorithms in all cases except that of a bi-directional heterogeneous ring, for which the problem remains open.

**Keywords:** Heterogeneous rings, data redistribution algorithms, load-balancing

## Résumé

Dans ce rapport, nous nous intéressons au problème de redistribution de données sur des anneaux de processeurs homogènes et hétérogènes. Ce problème surgit dans plusieurs applications, après chaque phase d'équilibrage de charge (nous ne discutons pas ici du mécanisme d'équilibrage de charge lui-même). Nous proposons des algorithmes qui visent à optimiser la redistribution de données pour des anneaux unidirectionnels et bidirectionnels, et nous donnons toutes les preuves de correction de ces algorithmes. Une des contributions principales de ce rapport est que nous pouvons prouver l'optimalité des algorithmes proposés dans tous les cas, sauf dans le cas d'un anneau hétérogène bidirectionnel, pour lequel le problème reste ouvert.

**Mots-clés:** Anneaux hétérogènes, algorithmes de redistribution de données, équilibrage de charge

# 1 Introduction

In this paper, we consider the problem of redistributing data on a heterogeneous ring of processors. The problem typically arises when a load balancing phase must be initiated. Because either of variations in the resource performances (CPU speed, communication bandwidth) or in the system/application requirements (completed tasks, new tasks, migrated tasks, etc.), data must be redistributed between participating processors so that the current (estimated) load is better balanced. We do not discuss the load-balancing mechanism itself (we take it as external, be it a system, an algorithm, an oracle or whatever). Rather we aim at optimizing the data redistribution induced by the load-balancing mechanism.

We adopt the following abstract view of the problem. There are  $n$  participating processors  $P_1, P_2, \dots, P_n$ . Each processor  $P_k$  initially holds  $L_k$  atomic data items. The load-balancing system/algorithm/oracle has decided that the new load of  $P_k$  should be  $L_k - \delta_k$ . If  $\delta_k > 0$ , this means that  $P_k$  now is overloaded and should send  $\delta_k$  data items to other processors; if  $\delta_k < 0$ ,  $P_k$  is under-loaded and should receive  $-\delta_k$  items from other processors. Of course there is a conservation law:  $\sum_{k=1}^n \delta_k = 0$ . The goal is to determine the required communications and to organize them (what we call the data redistribution) in minimal time.

We assume that the participating processors are arranged along a ring, either unidirectional or bidirectional, and either with homogeneous or heterogeneous link bandwidths, hence a total of four different frameworks to deal with. There are two main contexts in which processor rings are useful. The first context is those of many applications which operate on ordered data, and where the order needs to be preserved. Think of a large matrix whose columns are distributed among the processors, but with the condition that each processor operates on a slice of consecutive columns. An overloaded processor  $P_i$  can send its first columns to the processor  $P_j$  that is assigned the slice preceding its own slice (and  $P_j$  would append these columns to the end of its slice); similarly,  $P_i$  can send its last columns to the processor which is assigned the next slice; obviously, these are the only possibilities. In other words, the ordered uni-dimensional data distribution calls for a uni-dimensional arrangement of the processors, i.e., along a ring.

The second context that may call for a ring is the simplicity of the programming. Using a ring, either uni- or bi-directional, allows for a simpler management of the data to be redistributed. Data intervals can be maintained and updated to characterize each processor load. Finally, we observe that parallel machines with a rich but fixed interconnection topology (hypercubes, fat trees, grids, to quote a few) are on the decline. Heterogeneous cluster architectures, which we target in this paper, have a largely unknown interconnection graph, with includes gateways, backbones, and switches, and modeling the communication graph as a ring is a reasonable, if conservative, choice.

As stated above, we discuss four cases for the redistribution algorithms. We delay the formal statement of the redistribution problems until Section 2, but we summarize the main results as follows. In the simplest case, that of a unidirectional homogeneous ring, we derive an optimal algorithm, and we prove its correctness in full details. Because the target architecture is quite simple, we are able to provide explicit (analytical) formulas for the number of data sent/received by each processor. The same holds true for the case of a bidirectional homogeneous ring, but the algorithm becomes more complicated. When assuming heterogeneous communication links, we still derive an optimal algorithm for the unidirectional case, but we have to use an asynchronous formulation. However, we have to resort to heuristics based upon linear programming relaxation for the bidirectional case. We point out that one major contribution of the paper is the design of optimal algorithms, together with their formal proof of correctness: to the best of our knowledge, this is the first time that optimal algorithms are introduced.

The rest of the paper is organized as follows. In Section 2 we formally state the optimization problem. For homogeneous networks (all links have same capacity), the optimal algorithms are described in Section 3 (unidirectional ring) and in Section 5 (bidirectional ring). For heterogeneous networks, the optimal asynchronous unidirectional algorithm is presented in Section 4, and the linear-programming based optimal algorithm for *light* redistributions on bidirectional links is explained in Section 6. Section 7 is devoted to a survey of related work. In Section 8, we report some simulation results that confirm the usefulness of data redistributions. Finally, Section 9

concludes the paper and highlights future work directions.

## 2 Framework

We consider a set of  $n$  processors  $P_1, P_2, \dots, P_n$  arranged along a ring. The successor of  $P_i$  in the ring is  $P_{i+1}$ , and its predecessor is  $P_{i-1}$ , where all indices are taken modulo  $n$ . For  $1 \leq k, l \leq n$ ,  $C_{k,l}$  denotes the *slice* of consecutive processors  $C_{k,l} = P_k, P_{k+1}, \dots, P_{l-1}, P_l$ .

We denote by  $c_{i,i+1}$  the capacity of the communication link from  $P_i$  to  $P_{i+1}$ . In other words, it takes  $c_{i,i+1}$  time-units to send a data item from processor  $P_i$  to processor  $P_{i+1}$ . In the case of a bidirectional ring,  $c_{i,i-1}$  is the capacity of the link from  $P_i$  to  $P_{i-1}$ . We use the one-port model for communications: at any given time, there are at most two communications involving a given processor, one sent and the other received. A given processor can simultaneously send and receive data, so there is no restriction in the unidirectional case; however, in the bidirectional case, a given processor cannot simultaneously send data to its successor and its predecessor; neither can it receive data from both sides. These is the only restriction induced by the model: any pair of communications that does not violate the one-port constraint can take place in parallel.

Each processor  $P_k$  initially holds  $L_k$  atomic data items. After redistribution,  $P_k$  will hold  $L_k - \delta_k$  atomic data items. We call  $\delta_k$  the *unbalance* of  $P_k$ . We denote by  $\delta_{k,l}$  the total unbalance of the processor slice  $C_{k,l}$ :  $\delta_{k,l} = \delta_k + \delta_{k+1} + \dots + \delta_{l-1} + \delta_l$ .

Because of the conservation law of atomic data items,  $\sum_{k=1}^n \delta_k = 0$ . Obviously the unbalance cannot be larger than the initial load:  $L_k \geq \delta_k$ . In fact, we suppose that any processor holds at least one data, both initially ( $L_k \geq 1$ ) and after the redistribution ( $L_k \geq 1 + \delta_k$ ): otherwise we would have to build a new ring from the subset of resources still involved in the computation.

## 3 Homogeneous unidirectional ring

In this section, we consider a homogeneous unidirectional ring. Any processor  $P_i$  can only send data items to its successor  $P_{i+1}$ , and  $c_{i,i+1} = c$  for all  $i \in [1, n]$ . We first derive a lower bound on the running time of any redistribution algorithm. Then, we present an algorithm achieving this bound (hence optimal), and we prove its correctness.

### 3.1 Lower bound

We have the following bound on the optimal redistribution time:

**Lemma 1.** *Let  $\tau$  be the optimal redistribution time. Then:*

$$\tau \geq \left( \max_{1 \leq k \leq n, 0 \leq l \leq n-1} |\delta_{k,k+l}| \right) \times c. \quad (1)$$

*Proof.* The processor slice  $C_{k,k+l} = P_k, P_{k+1}, \dots, P_{k+l-1}, P_{k+l}$  has a total unbalance of  $\delta_{k,k+l} = \delta_k + \delta_{k+1} + \dots + \delta_{k+l-1} + \delta_{k+l}$ . If  $\delta_{k,k+l} > 0$ ,  $\delta_{k,k+l}$  data items must be sent from  $C_{k,k+l}$  to the other processors. The ring is unidirectional, so  $P_{k+l}$  is the only processor in  $C_{k,k+l}$  with an outgoing link. Furthermore,  $P_{k+l}$  needs a time equal to  $\delta_{k,k+l} \times c$  to send  $\delta_{k,k+l}$  data items. Therefore, in any case, a redistribution scheme cannot take less than  $\delta_{k,k+l} \times c$  to redistribute all data items. We have the same type of reasoning for the case  $\delta_{k,k+l} < 0$ .  $\square$

### 3.2 An optimal algorithm

We introduce the following redistribution algorithm:

We first prove the correction of Algorithm 1 (Lemma 3). Secondly, we prove its optimality (Lemma 4). Intuitively, if Step 6 of this algorithm is always feasible, then each execution of Step 3 has exactly a length of  $c$ , and the algorithm will meet the time bound of Lemma 1.

**Algorithm 1** Redistribution algorithm for homogeneous unidirectional rings

- 
- 1: Let  $\delta_{\max} = (\max_{1 \leq k \leq n, 0 \leq l \leq n-1} |\delta_{k,k+l}|)$
  - 2: Let **start** and **end** be two indices such that the slice  $C_{\text{start},\text{end}}$  is of maximal unbalance:  
 $\delta_{\text{start},\text{end}} = \delta_{\max}$ .
  - 3: **for**  $s = 1$  to  $\delta_{\max}$  **do**
  - 4:   **for all**  $l = 0$  to  $n - 1$  **do**
  - 5:     **if**  $\delta_{\text{start},\text{start}+l} \geq s$  **then**
  - 6:        $P_{\text{start}+l}$  sends to  $P_{\text{start}+l+1}$  a data item during the time interval  $[(s - 1) \times c, s \times c[$
- 

First, we point out that the slice  $C_{\text{start},\text{end}}$  is well-defined in Step 2 of the algorithm: for any slice with an unbalance  $\delta$ , the slice made up from the remaining processors has the opposite unbalance  $-\delta$ . Next, we state the particular role of the processor  $P_{\text{start}}$ :

**Lemma 2.** *Processor  $P_{\text{start}}$  receives no data items during the execution of Algorithm 1.*

*Proof.* We prove the result by contradiction. Suppose that at a given iteration  $s$  processor  $P_{\text{start}}$  receives some data items. Then the predecessor of  $P_{\text{start}}$  in the ring,  $P_{\text{start}-1}$ , sends a data item at this iteration. Thus,  $P_{\text{start}-1}$  being a sender, by the condition at Step 5 of Algorithm 1,  $\sum_{j=0}^{n-1} \delta_{\text{start}+j} = \delta_{\text{start},\text{start}-1} \geq s$ . However, due to the conservation law,  $\sum_{i=1}^n \delta_i = 0$ . Hence,  $0 \geq s$ , the desired contradiction.  $\square$

To prove that Algorithm 1 is correct, we must show that during each iteration, any processor required to send a data item in Step 6 actually holds at least one data item at this iteration. In other words, we must prove that no processor is asked to send a data item that it does not currently own. Let  $L_i^s$  be the load of  $P_i$  at the end of iteration  $s$  of Algorithm 1:

**Lemma 3.** *During iteration  $s$  of loop 3, if  $P_i$  sends a data item, then  $L_i^{s-1} \geq 1$ .*

*Proof.* We prove Lemma 3 by induction. Initially, by definition of unbalances (see Section 2), we know that each processor  $P_i$  in the ring initially holds an amount of  $L_i^0 = L_i \geq 1$  data items. Thus the result holds for  $s = 1$ .

Now we suppose that the result holds until a certain iteration  $s$  (included), and we focus on iteration  $s + 1$ . There are two cases to consider depending whether processor  $P_i$  is supposed to receive a data item during iteration  $s + 1$  or not:

1. If processor  $P_i$  is both a sender and a receiver during iteration  $s + 1$ , then  $P_i$  is both a sender and a receiver during iteration  $s$  by the condition at Step 5 of Algorithm 1. Then the load of  $P_i$  after iteration  $s$  was the same than before that iteration and  $L_i^s = L_i^{s-1}$ . We conclude using the induction hypothesis.
2. If processor  $P_i$  is a sender but not a receiver during iteration  $s + 1$ , we must verify that  $P_i$  does not send a data item that it does not own. Because  $P_i$  is a sender, then, by the condition at Step 5 of Algorithm 1, we have:

$$\delta_{\text{start},i} \geq s + 1. \quad (2)$$

Furthermore,  $P_i$  has sent a data item during each of the previous iterations.

During iteration  $s + 1$ ,  $P_i$  is not a receiver. Thus,  $P_{i-1}$  is not a sender during this iteration, and, by the condition at Step 5 of Algorithm 1, we have:  $\delta_{\text{start},i-1} < s + 1$ . During each iteration from 1 to  $\delta_{\text{start},i-1}$ ,  $P_{i-1}$  has sent a data item (see below for the proof that  $\delta_{\text{start},\text{start}+j} \geq 0$  for all  $j \in [0, n - 1]$ ). Hence, during each of these iterations,  $P_i$  was both a sender and a receiver, and neither its load nor its unbalance did change.

During each iteration from  $1 + \delta_{\text{start},i-1}$  to  $s$ , processor  $P_i$  was a sender but not a receiver. So both its load and its unbalance decrease by one during each of these iterations. Hence:

$$L_i^s = L_i - (s - \delta_{\text{start},i-1}). \quad (3)$$

However,  $\delta_i + \delta_{\text{start},i-1} = \delta_{\text{start},i}$ . So Equation 3 is equivalent to:  $L_i^s = L_i - \delta_i + \delta_{\text{start},i} - s$ . From Equation 2 we know that  $\delta_{\text{start},i} - s \geq 1$ . In Section 2, we assumed that  $L_i \geq 1 + \delta_i$ . So,  $L_i^s \geq 2$ .

The above proof relies on the property that, for any value of  $j \in [0, n-1]$ ,  $\delta_{\text{start},\text{start}+j} \geq 0$ . We now prove this result by contradiction. Hence we suppose that there exists a value  $j$  such that  $\delta_{\text{start},\text{start}+j} < 0$ . We have two cases to consider:

1.  $j + \text{start} \in [\text{start}, \text{end}]$ . Then  $\delta_{\text{start},\text{end}} = \delta_{\text{start},\text{start}+j} + \delta_{\text{start}+j+1,\text{end}}$  and  $\delta_{\text{start},\text{end}} < \delta_{\text{start}+j+1,\text{end}}$  which contradicts the maximality of  $C_{\text{start},\text{end}}$ .
2.  $j + \text{start} \notin [\text{start}, \text{end}]$ . Then  $\delta_{\text{start},j+\text{start}} = \delta_{\text{start},\text{end}} + \delta_{1+\text{end},j+\text{start}}$ . So  $\delta_{\text{start},\text{end}} < -\delta_{1+\text{end},j+\text{start}}$ . However, as the sum of unbalances is null by definition, the sum of unbalances of  $C_{1+\text{end},j+\text{start}}$  is equal to the opposite of the sum of unbalances of  $C_{j+1+\text{start},\text{end}}$ . Hence,  $\delta_{\text{start},\text{end}} < \delta_{j+1+\text{start},\text{end}}$ , which contradicts the maximality of  $C_{\text{start},\text{end}}$ .

□

We have proved the correction of Algorithm 1. We still have to prove that when it terminates, the entire redistribution has actually been performed:

**Lemma 4.** *When Algorithm 1 terminates after iteration  $\delta_{\max}$ , i.e., at time  $\tau$ , the load of any processor  $P_i$  is equal to  $L_i - \delta_i$ .*

*Proof.* We prove by induction on the processor indices, starting at processor  $P_{\text{start}}$ , that any processor  $P_j$  has the desired load of  $L_j - \delta_j$  at any iteration  $s \geq \max_{0 \leq i \leq j} \delta_{\text{start},\text{start}+i}$

As stated by Lemma 2, processor  $P_{\text{start}}$  never receives a data item during execution. So, after  $\delta_{\text{start},\text{start}} = \delta_{\text{start}}$  iterations of loop 3,  $P_{\text{start}}$  is never the receiver nor the sender of a data item. As required,  $P_{\text{start}}$  exactly holds  $L_{\text{start}} - \delta_{\text{start}}$  data items, i.e., its initial load minus the amount of data items sent.

We suppose the result proved up to a processor  $P_{\text{start}+l}$  (with  $l \geq 0$ ) included. We focus on processor  $P_{\text{start}+l+1}$ . Using the induction hypothesis, we know that at any iteration  $s \geq \max_{0 \leq i \leq l} \delta_{\text{start},\text{start}+i}$ , the total load of the slice  $C_{\text{start},\text{start}+l}$  is equal to  $\sum_{0 \leq i \leq l} L_i - \sum_{0 \leq i \leq l} \delta_i$ .

During the execution of the whole algorithm, processor  $P_{\text{start}+l+1}$  has sent exactly  $\delta_{\text{start},\text{start}+l+1}$  data items (remember that we showed in the proof of Lemma 3 that for any  $j \in [0, n-1]$ ,  $\delta_{\text{start},\text{start}+j} \geq 0$ ). All these send operations took place before or during iteration  $\delta_{\text{start},\text{start}+l+1}$ . Furthermore, Lemma 2 states that processor  $P_{\text{start}}$  never receives a data item during the execution. So, the total load of the slice  $C_{\text{start},\text{start}+l+1}$  does not change after iteration  $\delta_{\text{start},\text{start}+l+1}$ , and its total load is equal to its initial total load minus the data items sent by processor  $P_{\text{start}+l+1}$ :  $(\sum_{0 \leq i \leq l+1} L_i) - \delta_{\text{start},\text{start}+l+1}$ . Therefore, after any iteration  $s$ , where

$$s \geq \max \left( \max_{0 \leq i \leq l} \delta_{\text{start},\text{start}+i}, \delta_{\text{start},\text{start}+l+1} \right) = \max_{0 \leq i \leq l+1} \delta_{\text{start},\text{start}+i},$$

we know the total load of the slices  $C_{\text{start},\text{start}+l}$  and  $C_{\text{start},\text{start}+l+1}$ . Therefore, we know the load of processor  $P_{\text{start}+l+1}$ :

$$\begin{aligned} L_{\text{start}+l+1}^t &= \left( \left( \sum_{0 \leq i \leq l+1} L_{\text{start}+i} \right) - \delta_{\text{start},\text{start}+l+1} \right) - \left( \sum_{0 \leq i \leq l} L_{\text{start}+i} - \sum_{0 \leq i \leq l} \delta_{\text{start}+i} \right) \\ &= L_{\text{start}+l+1} - \delta_{\text{start}+l+1}. \end{aligned}$$

To conclude, we just need to remark that  $\delta_{\max} = \max_{0 \leq i \leq n-1} \delta_{\text{start},\text{start}+i}$ . □

The optimality of Algorithm 1 is a direct consequence of the previous lemmas:

**Theorem 1.** *Algorithm 1 is optimal.*

## 4 Heterogeneous unidirectional ring

In this section we still suppose that the ring is unidirectional but we no longer assume the communication paths to have the same capacities. We build on the results of the previous section to design an optimal algorithm (Algorithm 2 below). In this algorithm, the amount of data items sent by any processor  $P_i$  is exactly the same as in Algorithm 1 (namely  $\delta_{\text{start},i}$ ). However, as the communication links have different capabilities, we no longer have a synchronous behavior. A processor  $P_i$  sends its  $\delta_{\text{start},i}$  data items as soon as possible, but we cannot express its completion time with a simple formula. Indeed, if  $P_i$  initially holds more data items than it has to send, we have the same behavior than previously:  $P_i$  can send its data items during the time interval  $[0, \delta_{\text{start},i} \times c_{i,i+1}]$ . On the contrary, if  $P_i$  holds less data items than it has to send ( $L_i < \delta_{\text{start},i}$ ),  $P_i$  still starts to send some data items at time 0 but may have to wait to have received some other data items from  $P_{i-1}$  to be able to forward them to  $P_{i+1}$ .

---

### Algorithm 2 Redistribution algorithm for heterogeneous unidirectional rings

---

- 1: Let  $\delta_{\text{max}} = (\max_{1 \leq k \leq n, 0 \leq l \leq n-1} |\delta_{k,k+l}|)$
  - 2: Let **start** and **end** be two indices such that the slice  $C_{\text{start},\text{end}}$  is of maximal unbalance:  
 $\delta_{\text{start},\text{end}} = \delta_{\text{max}}$ .
  - 3: **for all**  $l = 0$  to  $n - 1$  **do**
  - 4:  $P_{\text{start}+l}$  sends  $\delta_{\text{start},\text{start}+l}$  data items one by one and as soon as possible to processor  $P_{\text{start}+l+1}$
- 

The asynchronousness of Algorithm 2 implies that it is correct by construction: we wait for receiving a data item before sending. Furthermore, when the algorithm terminates, the redistribution is complete (the proof is the same as in Lemma 4). There remains to prove that the running time of Algorithm 2 is optimal. We first compute this running time:

**Lemma 5.** *The running time of Algorithm 2 is  $\max_{0 \leq l \leq n-1} \delta_{\text{start},\text{start}+l} \times c_{\text{start}+l,\text{start}+l+1}$ .*

The result of Lemma 5 is surprising. Intuitively, it says that the running time of Algorithm 2 is equal to the maximum of the communication times of all the processors, if each of them initially stored locally all the data items it will have to send throughout the execution of the algorithm. In other words, there is no forwarding delay, whatever the initial distribution. The proof of Lemma 5 is technical and can be omitted at first reading.

*Proof.* We prove the result by contradiction, assuming that the running time of Algorithm 2, denoted as  $t_{\text{max}}$ , is strictly greater than  $\max_{0 \leq l \leq n-1} \delta_{\text{start},\text{start}+l} \times c_{\text{start}+l,\text{start}+l+1}$  (we assume that the algorithm starts running at time 0). Let  $P_i$  be any processor whose running time is  $t_{\text{max}}$ , i.e., let  $P_i$  be any processor which terminates the emission of its last data item at time  $t_{\text{max}}$ . By hypothesis,  $t_{\text{max}} > \delta_{\text{start},i} \times c_{i,i+1}$ . Therefore, there is some time during the running time of the algorithm at which processor  $P_i$  is not sending any data items to processor  $P_{i+1}$ . Let  $t_i$  denote the *latest* time at which  $P_i$  is not sending any data items. Then, by definition of  $t_i$ , from time  $t_i$  until the completion of the algorithm, processor  $P_i$  is continuously sending data items to  $P_{i+1}$ . Let  $n_i$  denote the number of data items that  $P_i$  sends during that interval. Note that we have  $t_{\text{max}} = t_i + n_i \times c_{i,i+1}$ . We now prove by induction that for any value of  $j \geq 1$ :

1. Processor  $P_{i-j}$  sends a data item to processor  $P_{i-j+1}$  during the time interval  $[t_i - \sum_{k=1}^j c_{i-k,i-k+1}, t_i - \sum_{k=1}^{j-1} c_{i-k,i-k+1}]$ .
2. Between time  $t_i - \sum_{k=1}^j c_{i-k,i-k+1}$  and the completion of the algorithm, processor  $P_{i-j}$  sends at least  $j + n_i$  data items to processor  $P_{i-j+1}$ .
3.  $c_{i-j,i-j+1} \leq c_{i,i+1}$ .
4. Right before time  $t_i - \sum_{k=1}^j c_{i-k,i-k+1}$ , processor  $P_{i-j}$  is not sending any data items to processor  $P_{i-j+1}$  (it is idle in sending).



Once we have proved these properties, the contradiction follows from considering processor  $P_{\text{start}}$ . Processor  $P_{\text{start}}$  only sends data items that it initially holds ( $\delta_{\text{start}} = \delta_{\text{start},\text{start}} \leq L_{\text{start}}$ ), and receives no data items from its predecessor in the ring. However, using the above properties, there is a value of  $j \geq 0$  such that  $\text{start} = i - j$ , and between time  $t_i - \sum_{k=1}^{j+1} c_{i-k,i-k+1}$  and the completion of the algorithm, processor  $P_{i-j-1}$  sends at least  $j + 1 + n_i$  data items to processor  $P_{i-j} = P_{\text{start}}$ . Hence the contradiction.

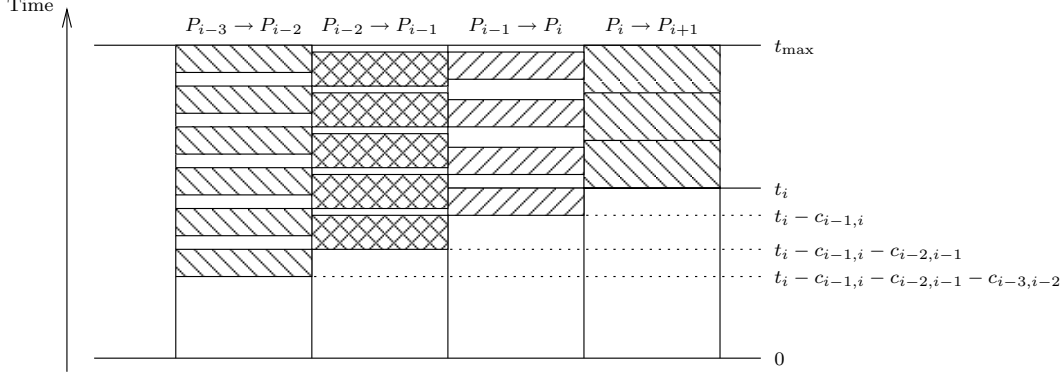


Figure 1: The construction used in the proof of Lemma 5.

The construction used in the proof is illustrated by Figure 1. We start by proving the above properties for  $j = 1$ .

1. By definition of  $t_i$ , processor  $P_i$  is not sending any data items to processor  $P_{i+1}$  right before time  $t_i$ . Because of the “as-soon-as” nature of the algorithm, processor  $P_i$  is not holding a single data item right before time  $t_i$  and is waiting for processor  $P_{i-1}$  to send it one. Furthermore, the data item that processor  $P_i$  started to send at time  $t_i$  is sent to it by processor  $P_{i-1}$  during the time interval  $[t_i - c_{i-1,i}, t_i]$ .
2. Between time  $t_i$  and the completion of the algorithm, processor  $P_i$  sends  $n_i$  data items to processor  $P_{i+1}$ . By hypothesis, processor  $P_i$  holds at least one data item after the completion of the algorithm. As  $P_i$  holds no data item right before time  $t_i$ , then between the times  $t_i - c_{i-1,i}$  and  $t_{\text{max}}$ ,  $P_{i-1}$  sends at least  $1 + n_i$  data items to  $P_i$ .
3. From what just precedes, and using the relationship between  $t_i$ ,  $n_i$ , and  $t_{\text{max}}$ , we infer:

$$t_i + n_i \times c_{i,i+1} = t_{\text{max}} \geq (t_i - c_{i-1,i}) + (1 + n_i) \times c_{i-1,i} \quad \Rightarrow \quad c_{i,i+1} \geq c_{i-1,i}$$

as  $n_i$  is nonzero by definition.

4. Suppose that processor  $P_{i-1}$  is sending a data item to processor  $P_i$  right before the time  $t_i - c_{i-1,i}$ . Then, at the earliest, this data item is received by processor  $P_i$  at time  $t_i - c_{i-1,i}$ . Due to the “as-soon-as” nature of the algorithm,  $P_i$  forwards this data item to processor  $P_{i+1}$  (as it forwards data items received later).  $P_i$  finishes to forward this data item at time  $t_i - c_{i-1,i} + c_{i,i+1} \geq t_i$  at the earliest. Therefore, processor  $P_i$  has no reason not to be sending any data item at time  $t_i$ , which contradicts the definition of  $t_i$ .

We now proceed to the general case of the induction. We suppose that the property is proved up to a processor  $P_{i-j}$  included (with  $j \geq 1$ ).

1. By induction hypothesis, processor  $P_{i-j}$  is not sending any data items to processor  $P_{i-j+1}$  right before time  $t_i - \sum_{k=1}^j c_{i-k,i-k+1}$ . Because of the “as-soon-as” nature of the algorithm, processor  $P_{i-j}$  is not holding a single data item right before this time and is waiting for processor  $P_{i-j-1}$  to send one. Furthermore, the data item that processor  $P_{i-j}$  started to send at time  $t_i - \sum_{k=1}^j c_{i-k,i-k+1}$  is sent to it by processor  $P_{i-j-1}$  during the time interval  $[t_i - \sum_{k=1}^{j+1} c_{i-k,i-k+1}, t_i - \sum_{k=1}^j c_{i-j,i-j+1}]$ .

2. Between time  $t_i - \sum_{k=1}^j c_{i-k, i-k+1}$  and the completion of the algorithm, processor  $P_{i-j}$  sends  $j + n_i$  data items to processor  $P_{i-j+1}$ , by induction hypothesis. By hypothesis, processor  $P_{i-j}$  holds at least one data item after the completion of the algorithm. As  $P_{i-j}$  holds no data item right before time  $t_i - \sum_{k=1}^j c_{i-k, i-k+1}$ , then between the times  $t_i - \sum_{k=1}^{j+1} c_{i-k, i-k+1}$  and  $t_{\max}$ ,  $P_{i-j-1}$  sends at least  $1 + j + n_i$  data items to  $P_{i-j}$ .
3. From what just precedes, and using the relationship between  $t_i$ ,  $n_i$ , and  $t_{\max}$ , we infer:

$$t_i + n_i \times c_{i, i+1} = t_{\max} \geq \left( t_i - \sum_{k=1}^{j+1} c_{i-k, i-k+1} \right) + (1 + j + n_i) \times c_{i-j-1, i-j} \quad \Rightarrow$$

$$n_i \times c_{i, i+1} + \sum_{k=1}^j c_{i-k, i-k+1} \geq (j + n_i) \times c_{i-j-1, i-j} \quad \Rightarrow$$

$$c_{i, i+1} \geq c_{i-j-1, i-j}$$

as, by induction hypothesis, for any  $k \in [1, j]$ ,  $c_{i, i+1} \geq c_{i-k, i-k+1}$ .

4. Suppose that processor  $P_{i-j-1}$  is sending a data item to processor  $P_{i-j}$  right before the time  $t_i - \sum_{k=1}^{j+1} c_{i-k, i-k+1}$ . Then, at the earliest, this data item is received by processor  $P_{i-j}$  at time  $t_i - \sum_{k=1}^{j+1} c_{i-k, i-k+1}$ . Due to the ‘‘as-soon-as’’ nature of the algorithm,  $P_{i-j}$  forwards this data item to processor  $P_{i-j+1}$  (as it forwards data items received later).  $P_{i-j}$  finishes to forward this data item at time  $t_i - c_{i-j-1, i-j} - \sum_{k=1}^{j-1} c_{i-k, i-k+1}$  at the earliest. Then, following the same line of reasoning, processor  $P_{i-j+1}$  forwards it to  $P_{i-j+2}$ , which receives it at the earliest at time  $t_i - c_{i-j-1, i-j} - \sum_{k=1}^{j-2} c_{i-k, i-k+1}$ , and so on. So, processor  $P_i$  receives this data item at the earliest at time  $t_i - c_{i-j-1, i-j}$ , and forwards it. Then, it finishes to send it at the earliest at time  $t_i - c_{i-j-1, i-j} + c_{i, i+1} \geq t_i$ , as we have seen that  $c_{i, i+1} \geq c_{i-j-1, i-j}$ . Therefore, processor  $P_i$  has no reason not to be sending any data items at time  $t_i$ , which contradicts the definition of  $t_i$ . Hence, processor  $P_{i-j-1}$  is not sending any data item to processor  $P_{i-j}$  right before the time  $t_i - \sum_{k=1}^{j+1} c_{i-k, i-k+1}$ .

□

**Theorem 2.** *Algorithm 2 is optimal.*

*Proof.* Let  $\tau$  denote the optimal redistribution time. Following the arguments used in the proof of Lemma 1 for the homogeneous case in Section 3.1, we obtain the lower bound:

$$\tau \geq \max_{1 \leq k \leq n, 0 \leq l \leq n-1} |\delta_{k, k+l}| \times c_{k+l, k+l+1}.$$

We conclude using Lemma 5. □

## 5 Homogeneous bidirectional ring

In this section, we consider a homogeneous bidirectional ring. All links have the same capacity but a processor can send data items to its two neighbors in the ring: there exists a constant  $c$  such that, for all  $i \in [1, n]$ ,  $c_{i, i+1} = c_{i, i-1} = c$ . We proceed as for the homogeneous unidirectional case: we first derive a lower bound on the running time of any redistribution algorithm, and then we present an algorithm attaining this bound.

## 5.1 Lower bound

We have the following bound on the optimal redistribution time:

**Lemma 6.** *Let  $\tau$  be the optimal redistribution time. Then:*

$$\tau \geq \max \left\{ \max_{1 \leq i \leq n} |\delta_i|, \max_{1 \leq i \leq n, 1 \leq l \leq n-1} \left\lceil \frac{|\delta_{i,i+l}|}{2} \right\rceil \right\} \times c. \quad (4)$$

*Proof.* Consider any processor  $P_i$  with positive unbalance ( $\delta_i > 0$ ). Even if processor  $P_i$  can send data items to both of its neighbors, because of the one-port model, it cannot send data items to both of them *simultaneously*. So, it requires processor  $P_i$  at least a time of  $\delta_i \times c$  to send  $\delta_i$  data items, whatever the destinations of these data items. We have a symmetric result for the case  $\delta_i < 0$ . Hence a first lower-bound on the optimal redistribution time  $\tau$ :

$$\tau \geq \left( \max_{1 \leq i \leq n} |\delta_i| \right) \times c.$$

Now, consider any non trivial slice of consecutive processors  $C_{k,l}$ . By “non trivial” we mean that the slice is not reduced to a single processor (we already treated that case) and that it does not contain all processors. We suppose that  $\delta_{k,l} > 0$ . So, in any redistribution scheme, at least  $\delta_{k,l}$  data items must be sent by  $C_{k,l}$ . As this slice is not reduced to a single processor, the two processors at the extremities of the slice,  $P_k$  and  $P_l$ , can simultaneously send data items to their neighbors outside of the slice,  $P_{k-1}$  and  $P_{l+1}$  respectively. Therefore, during any time interval of length  $c$ , at most two data items can be sent from the slice. So, it takes at least a time of  $\lceil \frac{\delta_{k,l}}{2} \rceil$  for the slice  $C_{k,l}$  to send  $\delta_{k,l}$  data items. Once again, the reasoning is similar when receiving data items if  $\delta_{k,l} < 0$ . Hence a second lower-bound on  $\tau$ :

$$\tau \geq \left( \max_{1 \leq i \leq n, 1 \leq l \leq n-1} \left\lceil \frac{|\delta_{i,i+l}|}{2} \right\rceil \right) \times c.$$

We just gather the previous two lower-bounds to obtain the desired bound.  $\square$

## 5.2 An optimal algorithm

Algorithm 3 (see below) is a recursive algorithm which defines communication patterns designed so as to decrease the value of  $\delta_{\max}$  (computed at Step 1) by one from one recursive call to another. The intuition behind Algorithm 3 is the following:

1. Any non trivial slice  $C_{k,l}$  such that  $\lceil \frac{|\delta_{k,l}|}{2} \rceil = \delta_{\max}$  and  $\delta_{k,l} \geq 0$  must send two data items per recursive call, one through each of its extremities.
2. Any non trivial slice  $C_{k,l}$  such that  $\lceil \frac{|\delta_{k,l}|}{2} \rceil = \delta_{\max}$  and  $\delta_{k,l} \leq 0$  must receive two data items per recursive call, one through each of its extremities.
3. Once the mandatory communications specified by the two previous cases are defined, we take care of any processor  $P_i$  such that  $|\delta_i| = \delta_{\max}$ . If  $P_i$  is already involved in a communication due to the previous cases, everything is settled. Otherwise, we have the freedom to choose whom  $P_i$  will send a data item to (case  $\delta_i > 0$ ) or whom  $P_i$  will receive a data item from (case  $\delta_i < 0$ ). To simplify the algorithm we decide that all these communications will take place in the direction from  $P_i$  to  $P_{i+1}$ .

Algorithm 3 is initially called with the parameter  $s = 1$ . For any call to Algorithm 3, all the communications take place in parallel and exactly at the same time, because the communication paths are homogeneous by hypothesis. One very important point about Algorithm 3 is that this algorithm is a set of rules which *only* specify which processor  $P_i$  must send a data item to which processor  $P_j$ , one of its immediate neighbors. Therefore, whatever the number of rules deciding

**Algorithm 3** Redistribution algorithm for homogeneous bidirectional rings (for step  $s$ )

---

```

1: Let  $\delta_{\max} = \max\{\max_{1 \leq i \leq n} |\delta_i|, \max_{1 \leq i \leq n, 1 \leq l \leq n-1} \lceil \frac{|\delta_{i,i+l}|}{2} \rceil\}$ 
2: if  $\delta_{\max} \geq 1$  then
3:   if  $\delta_{\max} \neq 2$  then
4:     for all slice  $C_{k,l}$  such that  $\delta_{k,l} > 1$  and  $\lceil \frac{|\delta_{k,l}|}{2} \rceil = \delta_{\max}$  do
5:        $P_k$  sends a data item to  $P_{k-1}$  during the time interval  $[(s-1) \times c, s \times c[$ .
6:        $P_l$  sends a data item to  $P_{l+1}$  during the time interval  $[(s-1) \times c, s \times c[$ .
7:     for all slice  $C_{k,l}$  such that  $\delta_{k,l} < -1$  and  $\lceil \frac{|\delta_{k,l}|}{2} \rceil = \delta_{\max}$  do
8:        $P_{k-1}$  sends a data item to  $P_k$  during the time interval  $[(s-1) \times c, s \times c[$ .
9:        $P_{l+1}$  sends a data item to  $P_l$  during the time interval  $[(s-1) \times c, s \times c[$ .
10:   else if  $\delta_{\max} = 2$  then
11:     for all slice  $C_{k,l}$  such that  $\delta_{k,l} \geq 3$  do
12:        $P_l$  sends a data item to  $P_{l+1}$  during the time interval  $[(s-1) \times c, s \times c[$ .
13:     for all slice  $C_{k,l}$  such that  $\delta_{k,l} = 4$  do
14:        $P_k$  sends a data item to  $P_{k-1}$  during the time interval  $[(s-1) \times c, s \times c[$ .
15:     for all slice  $C_{k,l}$  such that  $\delta_{k,l} \leq -3$  do
16:        $P_{k-1}$  sends a data item to  $P_k$  during the time interval  $[(s-1) \times c, s \times c[$ .
17:     for all slice  $C_{k,l}$  such that  $\delta_{k,l} = -4$  do
18:        $P_{l+1}$  sends a data item to  $P_l$  during the time interval  $[(s-1) \times c, s \times c[$ .
19:   for all processor  $P_i$  such that  $\delta_i = \delta_{\max}$  do
20:     if  $P_i$  is not already sending, due to one of the previous steps, a data item during the time interval  $[(s-1) \times c, s \times c[$  then
21:        $P_i$  sends a data item to  $P_{i+1}$  during the time interval  $[(s-1) \times c, s \times c[$ .
22:   for all processor  $P_i$  such that  $\delta_i = -(\delta_{\max})$  do
23:     if  $P_i$  is not already receiving, due to one of the previous steps, a data item during the time interval  $[(s-1) \times c, s \times c[$  then
24:        $P_i$  receives a data item from  $P_{i-1}$  during the time interval  $[(s-1) \times c, s \times c[$ .
25:   if  $\delta_{\max} = 1$  then
26:     for all processor  $P_i$  such that  $\delta_i = 0$  do
27:       if  $P_{i-1}$  sends a data item to  $P_i$  during the time interval  $[(s-1) \times c, s \times c[$  then
28:          $P_i$  sends a data item to  $P_{i+1}$  during the time interval  $[(s-1) \times c, s \times c[$ .
29:       if  $P_{i+1}$  sends a data item to  $P_i$  during the time interval  $[(s-1) \times c, s \times c[$  then
30:          $P_i$  sends a data item to  $P_{i-1}$  during the time interval  $[(s-1) \times c, s \times c[$ .
31:   Recursive call to Algorithm 3 ( $s+1$ )

```

---

that there must be some data item sent from a processor  $P_i$  to one of its immediate neighbor  $P_j$ , only one data item is sent from  $P_i$  to  $P_j$  to satisfy all these rules.

To prove that Algorithm 3 is optimal, we show that the set of rules is consistent, i.e., that it respects the one-port model, and that the value  $\delta_{\max}$  (computed at Step 1) decreases by one at each recursive call.

**Lemma 7.** *Algorithm 3 satisfies to all the one-port constraints.*

*Proof.* We call *maximal slice* a slice  $C_{k,l}$  of consecutive processors whose total unbalance satisfies the condition:  $\lceil \frac{|\delta_{k,l}|}{2} \rceil = \delta_{\max}$ . We call *maximal processor* a processor  $P_i$  whose unbalance is equal to  $\delta_{\max}$  or  $-\delta_{\max}$ :  $|\delta_i| = \delta_{\max}$ . Maximal slices are processed by rules at Steps 4 through 18, while maximal processors are processed by the rules of Steps 19 and 22.

To prove that the set of rules obeys the one-port model, we have to prove that no processor simultaneously receives one data item from both neighbors, and that no processor simultaneously sends one data item to both neighbors. We only study the cases involving a processor receiving data items from both neighbors, because the algorithm symmetrically processes sends and receives.

We prove the result by contradiction. So, suppose that there exists a processor  $P_j$  that receives one data item from each neighbor,  $P_{j-1}$  and  $P_{j+1}$ . There are four cases to consider:

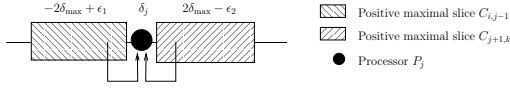


Figure 2: Case 1a in the proof of Lemma 7.

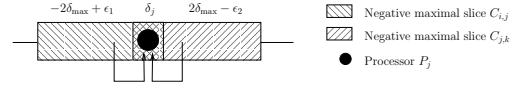


Figure 3: Case 1b in the proof of Lemma 7.

1.  $P_{j-1}$  and  $P_{j+1}$  are both sending a data item to  $P_j$  because of Steps 4 through 18. Then,  $P_{j-1}$  and  $P_{j+1}$  send data items to  $P_j$  either because they are extremities of positive maximal slices or because  $P_j$  is the extremity of (a) negative maximal slice(s). We thus have three subcases to study:

- (a)  $P_{j-1}$  and  $P_{j+1}$  are both extremities of positive maximal slices. Then there exist two indices  $i$  and  $k$  such that the slices  $C_{i,j-1}$  and  $C_{j+1,k}$  are both positive maximal slices. So, by definition, there exist two values  $\epsilon_1$  and  $\epsilon_2$ , each one either equal to 0 or 1, such that  $\delta_{i,j-1} = 2\delta_{\max} - \epsilon_1$  and  $\delta_{j+1,k} = 2\delta_{\max} - \epsilon_2$ . This case is illustrated by Figure 2. Consider the slice  $C_{i,k}$ . By definition of  $\delta_{\max}$  we have:

$$\begin{aligned} \left\lceil \frac{\delta_{i,k}}{2} \right\rceil &\leq \delta_{\max} && \Leftrightarrow \\ \left\lceil \frac{(2\delta_{\max} - \epsilon_1) + \delta_j + (2\delta_{\max} - \epsilon_2)}{2} \right\rceil &\leq \delta_{\max} && \Leftrightarrow \\ 4\delta_{\max} + \delta_j - \epsilon_1 - \epsilon_2 &\leq 2\delta_{\max} && \Leftrightarrow \\ \delta_j &\leq \epsilon_1 + \epsilon_2 - 2\delta_{\max} \end{aligned}$$

However, by definition of  $\delta_{\max}$ ,  $\delta_j$  is greater than or equal to  $-\delta_{\max}$ . So we end up with the constraint:

$$-\delta_{\max} \leq \epsilon_1 + \epsilon_2 - 2\delta_{\max} \quad \Leftrightarrow \quad \delta_{\max} \leq \epsilon_1 + \epsilon_2. \quad (5)$$

We then have three cases two consider:

- i.  $\delta_{\max} = 0$ : there is nothing to do as stated by the test at Step 2. (In the remaining of this proof, we will no more consider the cases where  $\delta_{\max} = 0$ .)
  - ii.  $\delta_{\max} = 1$ . Then, either  $\epsilon_1 = 1$  and  $\delta_{i,j-1} = 1$ , or  $\epsilon_2 = 1$  and  $\delta_{j+1,k} = 1$ : in both cases, this contradicts our hypothesis that  $P_{j-1}$  and  $P_{j+1}$  are both sending data items to  $P_j$  because of Steps 4 through 18.
  - iii.  $\delta_{\max} = 2$ . This case is illustrated by Figure 4. Equation 5 induces that  $\epsilon_1 = \epsilon_2 = 1$ . Applying the general scheme defined by Steps 4 through 6 would lead to the violation of the one-port model (cf. Figure 4(a)). However, each of the two slices  $C_{i,j-1}$  and  $C_{j+1,k}$  only needs to output three data items in two successive calls to Algorithm 3 (the calls with  $\delta_{\max} = 2$  and  $\delta_{\max} = 1$ ). So, we only require these maximal slices to output one data item during the call with  $\delta_{\max} = 2$ , in the direction from  $P_i$  to  $P_{i+1}$  (cf. Figures 4(b) and 4(c)). Remark that, in our example  $P_i$  outputs a data item at step  $\delta_{\max} = 2$ : this is not because it is the extremity of  $C_{i,j-1}$  with  $\delta_{i,j-1} = 3$  but because  $\delta_{i,k} = 4$ . This particular case is one of the reasons why we introduced the special processing of Steps 10 through 18.
- (b)  $P_j$  is the extremity of two negative maximal slices  $C_{i,j}$  and  $C_{j,k}$  with  $i < j < k$ . So, by definition, there exist two values  $\epsilon_1$  and  $\epsilon_2$ , each one either equal to 0 or 1, such that  $\delta_{i,j} = -2\delta_{\max} + \epsilon_1$  and  $\delta_{j,k} = -2\delta_{\max} + \epsilon_2$ . This case is illustrated by Figure 3. Consider the slice  $C_{i,k}$ :

$$\delta_{i,k} = \delta_{i,j} + \delta_{j,k} - \delta_j = -4\delta_{\max} + \epsilon_1 + \epsilon_2 - \delta_j$$

By definition of  $\delta_{\max}$  we have:  $\delta_{i,k} \geq -2\delta_{\max}$ . So,  $\delta_j \leq -2\delta_{\max} + \epsilon_1 + \epsilon_2$ . But  $\delta_j \geq -\delta_{\max}$ . Hence,  $\delta_{\max} \leq \epsilon_1 + \epsilon_2$ . We then have two cases two consider:

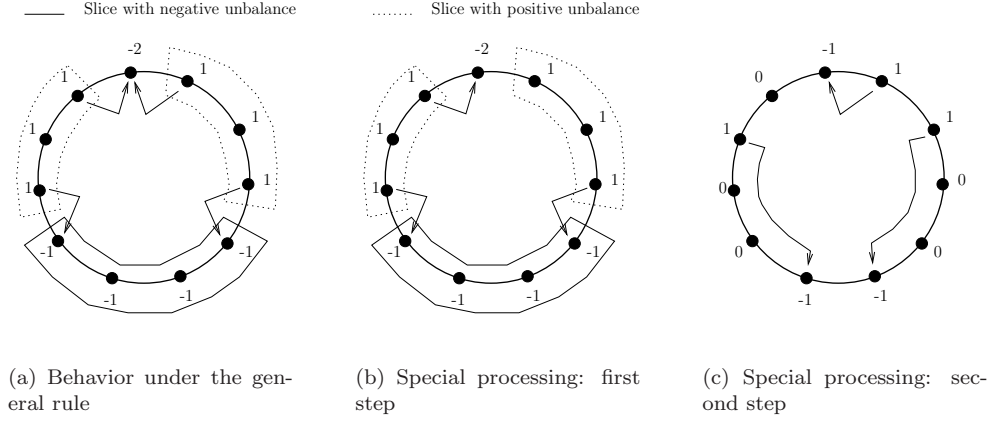


Figure 4: Case 1(a)iii in the proof of Lemma 7. Figure 4(a) shows the problem: the one-port model is violated if we apply the general rules to that case. Figures 4(b) and 4(c) describes the two steps of the special processing: in the first step, only one data item is output by the rightmost maximal slice; and in the second step, only one data item is output by the slice which was the leftmost maximal slice.

- i.  $\delta_{\max} = 1$ . Then, either  $\epsilon_1 = 1$  and  $\delta_{i,j} = -1$ , or  $\epsilon_2 = 1$  and  $\delta_{j,k} = -1$ . In both cases, this contradicts our hypothesis that  $P_{j-1}$  and  $P_{j+1}$  are both sending data items to  $P_j$  because of Steps 4 through 18.
- ii.  $\delta_{\max} = 2$ . Then  $\epsilon_1 = \epsilon_2 = 1$  and  $\delta_{i,j} = \delta_{j,k} = -3$ . As  $\delta_j \leq -2\delta_{\max} + \epsilon_1 + \epsilon_2$  and as, by definition of  $\delta_{\max}$ ,  $\delta_j \geq -\delta_{\max}$ , then  $\delta_j = -\delta_{\max} = -2$ .

Applying the general scheme defined by Steps 7 through 9 would lead to the violation of the one-port model (see Figure 5(a)). However, each of the two slices  $C_{i,j}$  and  $C_{j,k}$  only needs to input three data items in two successive calls to Algorithm 3 (the calls with  $\delta_{\max} = 2$  and  $\delta_{\max} = 1$ ). So, we only require these maximal slices to input one data item during the call with  $\delta_{\max} = 2$ , in the direction from  $P_i$  to  $P_{i+1}$  (see Figures 5(b) and 5(c)). Remark that, in our example  $P_k$  inputs a data item at step  $\delta_{\max} = 2$ : this is not because it is the extremity of  $C_{j,k}$  with  $\delta_{j,k} = -3$  but because  $\delta_{i,k} = -4$ .

This particular case is one of the reasons why we introduced the special processing of Steps 10 through 18.

- (c)  $P_j$  is the extremity of a negative maximal slice and one of its neighbor is the extremity of a positive maximal slice. Without any loss of generality, suppose that  $P_{j+1}$  sends a data item to  $P_j$  because  $C_{i,j}$  is a maximal negative slice. Then  $P_{j-1}$  sends a data item to  $P_j$  because it is the extremity of some positive maximal chain  $C_{k,j-1}$ . So, by definition, there exist two values  $\epsilon_1$  and  $\epsilon_2$ , each one either equal to 0 or 1, such that  $\delta_{i,j} = -2\delta_{\max} + \epsilon_1$  and  $\delta_{k,j-1} = 2\delta_{\max} - \epsilon_2$ . We have two cases to consider, depending whether the slice  $C_{k,j-1}$  is enclosed in the slice  $C_{i,j}$ :

- i.  $k \in [i, j-2]$  (this case is illustrated by Figure 6).  $\delta_{i,k-1} + \delta_j = \delta_{i,j} - \delta_{k,j-1} = (-2\delta_{\max} + \epsilon_1) - (2\delta_{\max} - \epsilon_2) = -4\delta_{\max} + \epsilon_1 + \epsilon_2$ . However, by definition of  $\delta_{\max}$ ,  $\delta_{i,k-1} \geq -2\delta_{\max}$  and  $\delta_j \geq -\delta_{\max}$ . So,  $\delta_{\max} \leq \epsilon_1 + \epsilon_2$ . Once again, we have two cases to consider:
  - A.  $\delta_{\max} = 1$ . Then, as always, either  $\epsilon_1 = 1$  and  $\delta_{i,j} = -1$ , or  $\epsilon_2 = 1$  and  $\delta_{j,k} = 1$ . In both cases, this contradicts our hypotheses on  $C_{i,j}$  and  $C_{k,j-1}$ .
  - B.  $\delta_{\max} = 2$ . Then  $\epsilon_1 = \epsilon_2 = 1$  and  $\delta_{i,j} = -3$  and  $\delta_{k,j-1} = 3$ . Therefore,  $\delta_{i,k-1} + \delta_j = -6$ . By definition of  $\delta_{\max}$ ,  $\delta_j \geq -\delta_{\max} = -2$  and  $\delta_{i,k-1} \geq -2\delta_{\max} = -4$ , we have  $\delta_j = -2$  and  $\delta_{i,k-1} = -4$ .

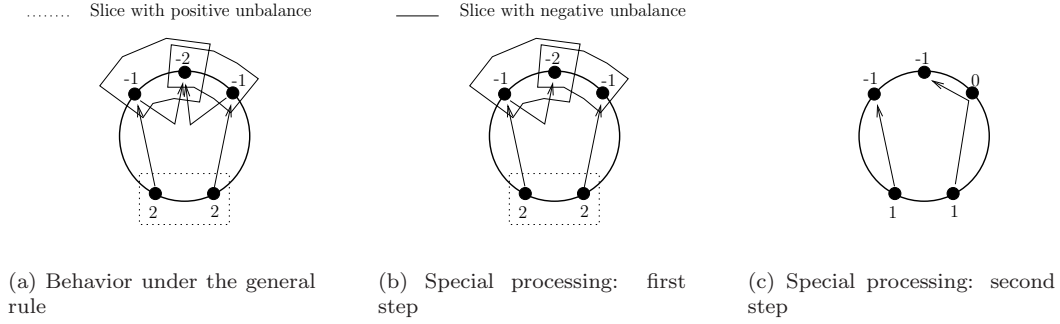


Figure 5: Case 1(b)ii in the proof of Lemma 7. Figure 5(a) shows the problem: the one-port model is violated if we apply the general rules to that case. Figures 5(b) and 5(c) describes the two steps of the special processing: in the first step, only one data item is input by the leftmost negative maximal slice; and in the second step, only one data item is input by the slice which was the rightmost maximal slice.

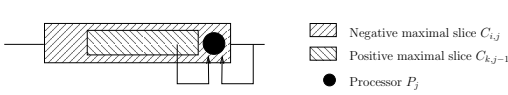


Figure 6: Case 1(c)i of the proof of Lemma 7.

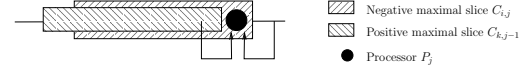


Figure 7: Case 1(c)ii of the proof of Lemma 7.

Similarly to the cases 1(a)iii and 1(b)ii, applying the general scheme defined by Steps 4 through 9 would lead to the violation of the one-port model (cf. Figure 8(a)). However, the slice  $C_{i,j}$  only needs to input three data items in two successive calls to Algorithm 3 (the calls with  $\delta_{\max} = 2$  and  $\delta_{\max} = 1$ ) while the slice  $C_{k,j-1}$  only needs to output three data items. So, we only require the slice  $C_{i,j}$  to input one data item and the slice  $C_{k,j-1}$  to output one data item during the call with  $\delta_{\max} = 2$ , both communications being in the direction from  $P_i$  to  $P_{i+1}$  (cf. Figures 8(b) and 8(c)). Remark that, in our example,  $P_k$  outputs a data item at step  $\delta_{\max} = 2$ : this is not because it is the extremity of  $C_{k,j-1}$  with  $\delta_{k,j-1} = 3$  but because  $\delta_{i,k-1} = -4$ .

This particular case is one of the reasons why we introduced the special processing of Steps 10 through 18.

- ii.  $k < i$  (this case is illustrated by Figure 7). Then,  $\delta_{k,i-1} = \delta_j + \delta_{k,j-1} - \delta_{i,j} = \delta_j + (2\delta_{\max} - \epsilon_1) - (2\delta_{\max} + \epsilon_2) = \delta_j + 4\delta_{\max} - \epsilon_1 - \epsilon_2$ . However, by definition of  $\delta_{\max}$ ,  $\delta_{k,i-1} \leq 2\delta_{\max}$  and  $\delta_j \geq -\delta_{\max}$ . So,  $\delta_j \leq -2\delta_{\max} - \epsilon_1 - \epsilon_2$  and, thus,  $-\delta_{\max} \leq -2\delta_{\max} - \epsilon_1 - \epsilon_2$ . Hence,  $\delta_{\max} = \epsilon_1 = \epsilon_2 = 0$ , which is absurd.

2.  $P_{j-1}$  and  $P_{j+1}$  are both sending data items to  $P_j$ : one sends data items due to Steps 4 through 18; the other one is a maximal processor which sends data items due to Steps 19 and 24. Without loss of generality, suppose that  $P_{j-1}$  is the maximal processor.

We have two cases to consider, depending whether  $P_{j+1}$  is sending a data item to  $P_j$  because of a positive or negative maximal slice.

- (a)  $P_{j+1}$  is the extremity of a positive maximal slice  $C_{j+1,k}$ . Figure 9 illustrates this case. Therefore, there exists  $\epsilon \in \{0; 1\}$ , such that  $\delta_{j+1,k} = 2\delta_{\max} - \epsilon$ . By hypothesis,  $P_{j-1}$  sends data items due to Steps 19 and 24, and thus the slice  $C_{j-1,k}$  is not a maximal

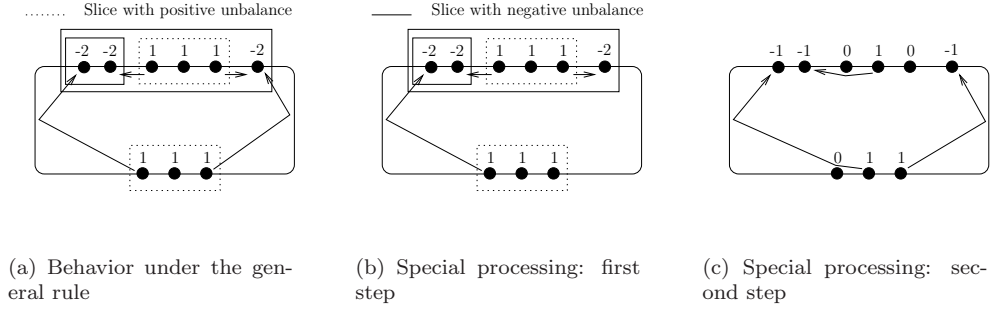


Figure 8: Case 1(c)iB in the proof of Lemma 7. Figure 8(a) shows the problem: the one-port model is violated if we apply the general rules to that case. Figures 8(b) and 8(c) describes the two steps of the special processing: in the first step, only one data item is input by the negative maximal slice; and in the second step, only one data item is output by the slice which was the positive maximal slice.

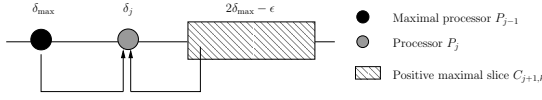


Figure 9: Case 2a of the proof of Lemma 7.

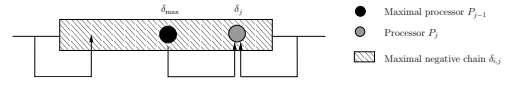


Figure 10: Case 2b of the proof of Lemma 7.

slice, i.e.,  $\lceil \frac{\delta_{j-1,k}}{2} \rceil \leq \delta_{\max} - 1$ .

$$\begin{aligned} \lceil \frac{\delta_{j-1,k}}{2} \rceil &= \left\lceil \frac{\delta_{\max} + \delta_j + 2\delta_{\max} - \epsilon}{2} \right\rceil \leq \delta_{\max} - 1 \\ \Leftrightarrow \delta_{\max} + \delta_j + 2\delta_{\max} - \epsilon &\leq 2\delta_{\max} - 2 \\ \Leftrightarrow \delta_j &\leq \epsilon - 2 - \delta_{\max} \\ \Rightarrow \delta_j &\leq -1 - \delta_{\max} \end{aligned}$$

which contradicts the definition of  $\delta_{\max}$ .

- (b)  $P_j$  is the extremity of a negative maximal slice  $C_{i,j}$  (Figure 10 illustrates this case). Then, there exists  $\epsilon \in \{0; 1\}$ , such that  $\delta_{i,j} = -2\delta_{\max} + \epsilon$ . Therefore,  $\delta_{i,j-2} = \delta_{i,j} - \delta_{\max} - \delta_j = -3\delta_{\max} + \epsilon - \delta_j$ . By hypothesis,  $P_{j-1}$  sends data items due to Steps 19 and 24, and thus the slice  $C_{i,j-2}$  is not a maximal slice, which implies that  $\delta_{i,j-2} \geq -2\delta_{\max} + 2$ . Therefore,  $-3\delta_{\max} + \epsilon - \delta_j \geq -2\delta_{\max} + 2$  and thus  $\delta_j \leq -\delta_{\max} + \epsilon - 2 \leq -\delta_{\max} - 1$ , which contradicts the definition of  $\delta_{\max}$ .
3.  $P_{j-1}$  and  $P_{j+1}$  are both sending data items to  $P_j$  because both are maximal processors which send data items due to Steps 19 through 24. This case is impossible as these steps only define data item sending in the direction from  $P_i$  to  $P_{i+1}$  and never in the reverse direction (from  $P_i$  to  $P_{i-1}$ ).
  4.  $P_j$  is a maximal processor of negative unbalance and this is the reason why  $P_{j-1}$  sends it a data item (following Steps 19 through 24). There maybe several reasons why  $P_{j+1}$  would also send a data item to  $P_j$ :
    - (a)  $P_{j+1}$  is the extremity of a positive maximal slice  $C_{j+1,k}$  and it sends a data item due to Steps 4 through 18. Then the test at Step 23 contradicts our hypothesis on  $P_{j-1}$ .
    - (b)  $P_{j+1}$  is a positive maximal processor. But in this case  $P_{j+1}$  sends a data item to  $P_{j+2}$  and not to  $P_j$ .



- (c)  $P_j$  is the extremity of a negative maximal slice  $C_{i,j}$  and  $P_{j+1}$  sends it a data item due to Steps 4 through 18. Then the test at Step 23 contradicts our hypothesis on  $P_j$ .  $\square$

**Lemma 8.** *Algorithm 3 terminates in exactly  $\max \left\{ \max_{1 \leq i \leq n} |\delta_i|, \max_{1 \leq i \leq n, 1 \leq l \leq n-1} \left\lceil \frac{\delta_{i,i+l}}{2} \right\rceil \right\}$  recursive calls.*

*Proof.* We prove that from one recursive call to Algorithm 3 to another, the value of  $\delta_{\max}$  (computed at Step 1) decreases by one. Therefore, we consider how unbalances change between the initial call to Algorithm 3 and its recursive call (excluded). For the general case, we have to prove four properties:

1. If the non-trivial slice  $C_{k,l}$  was initially a maximal slice, i.e., if  $\lceil \frac{|\delta_{k,l}|}{2} \rceil = \delta_{\max}$ , then after the communications we have  $\lceil \frac{|\delta_{k,l}|}{2} \rceil = \delta_{\max} - 1$ .

As previously we focus on the case of a positive maximal slice. The rules of Algorithm 3 are written so that the slice  $C_{k,l}$  sends two data items (or only one in the degenerate case when  $\delta_{\max} = 2$  and  $\delta_{k,l} = 3$ ) during an execution of Algorithm 3. This is all we need to conclude, provided that this slice does not receive any data item during this call.

Thus, suppose that  $C_{k,l}$  receives a data item. We have three cases to consider:

- (a) The maximal slice  $C_{k,l}$  receives a data item from a processor which is the extremity of another maximal slice and which sends a data item due to Steps 4 through 18. As the other maximal slice is sending a data item, its unbalance is positive. Without any loss of generality, we suppose it is a maximal slice of the form  $C_{l+1,m}$ . Then, by definition of maximal slices,  $\delta_{k,l} = 2\delta_{\max} - \epsilon_1$  and  $\delta_{l+1,m} = 2\delta_{\max} - \epsilon_2$ , with both  $\epsilon_1$  and  $\epsilon_2$  taking values in  $\{0, 1\}$ . Thus,  $\delta_{k,m} = 4\delta_{\max} - \epsilon_1 - \epsilon_2$ . However, by definition of  $\delta_{\max}$ ,  $\delta_{k,m} \leq 2\delta_{\max}$ . Hence, we obtain  $2\delta_{\max} \leq \epsilon_1 + \epsilon_2$ , which implies  $\delta_{\max} = 1$  and  $\epsilon_1 = \epsilon_2 = 1$ . Then,  $\delta_{l+1,m} = 1$  which contradicts the hypothesis that  $C_{l+1,m}$  sends a data item due to Steps 4 through 18 (see the test at Step 4).
- (b) The maximal slice  $C_{k,l}$  receives a data item from a processor which is maximal and which sends data items because of Steps 19 through 24. This case can only arise if this maximal processor has a positive unbalance. Without any loss of generality, we suppose processor  $P_{k-1}$  has an unbalance of  $\delta_{\max}$ . Then, by definition of maximal slices,  $\delta_{k,l} \geq 2\delta_{\max} - 1$  and  $\delta_{k-1,l} \geq 3\delta_{\max} - 1$ . However, by definition of  $\delta_{\max}$ ,  $\delta_{k-1,l} \leq 2\delta_{\max}$ . So  $\delta_{\max} = 1$  and  $\delta_{k-1,l} = 1$ . Then we have  $\delta_{k,l} = 1$ , and  $\delta_{k-1,l} = 2$ . Therefore,  $C_{k-1,l}$  is a maximal slice and processor  $P_{k-1}$  sends a data item to processor  $P_{k-2}$  rather than to  $P_k$ .
- (c) The maximal slice  $C_{k,l}$  receives a data item because one of its extremities is also the extremity of a negative maximal slice. Without any loss of generality, we suppose this negative maximal slice is of the form  $C_{k,m}$  (with  $l \in [k; m]$ ).  $C_{k,l}$  being a positive maximal slice,  $\delta_{k,l} = 2\delta_{\max} - \epsilon_1$  with  $\epsilon_1 \in \{0; 1\}$ .  $C_{k,m}$  being a negative maximal slice,  $\delta_{k,m} = -2\delta_{\max} + \epsilon_2$  with  $\epsilon_2 \in \{0; 1\}$ . We have two cases to consider:
  - i.  $l < m$ . Then,  $\delta_{l+1,m} = (-2\delta_{\max} + \epsilon_2) - (2\delta_{\max} - \epsilon_1) = -4\delta_{\max} + \epsilon_1 + \epsilon_2$ . By definition of  $\delta_{\max}$ ,  $\delta_{l+1,m} \geq -2\delta_{\max}$ , and thus  $\delta_{\max} = 1$  and  $\epsilon_1 = \epsilon_2 = 1$ . Then  $C_{k,m} = -1$  and, because of the test of Step 7, the rules of Steps 8 and 9 do not apply, and the maximal slice  $C_{k,l}$  does not receive a data item because it is enclosed in a negative maximal slice.
  - ii.  $l > m$ . Then,  $\delta_{m+1,l} = (2\delta_{\max} - \epsilon_1) - (-2\delta_{\max} + \epsilon_2) = 4\delta_{\max} - \epsilon_1 - \epsilon_2$ . By definition of  $\delta_{\max}$ ,  $\delta_{m+1,l} \leq 2\delta_{\max}$ , and thus  $\delta_{\max} = 1$  and  $\epsilon_1 = \epsilon_2 = 1$ . Then  $C_{k,m} = -1$ .

Thus, in both cases,  $C_{k,m} = -1$ . Then, because of the test of Step 7, the rules of Steps 8 and 9 do not apply, and the maximal slice  $C_{k,l}$  does not receive a data item because one of its extremities is also the extremity of a negative maximal slice.

2. If processor  $P_i$  was initially maximal, i.e., if  $|\delta_i| = \delta_{\max}$ , then after the communications we have  $|\delta_i| = \delta_{\max} - 1$ .

As previously, we only focus on the case  $\delta_i = \delta_{\max}$ . If, after communications, we do not have  $|\delta_i| = \delta_{\max} - 1$ , then  $P_i$  has received one data item.

- (a)  $P_i$  receives a data item from a processor which is the extremity of a positive maximal slice and which sends a data item due to Steps 4 through 18. Without loss of generality, suppose this processor is  $P_{i+1}$  and the slice  $C_{i+1,j}$ . By definition of maximal slices, there exists a value  $\epsilon$ , either equal to 0 or 1, such that  $\delta_{i+1,j} = 2\delta_{\max} - \epsilon$ . Then  $\delta_{i,j} = 3\delta_{\max} - \epsilon$ . As, by definition of  $\delta_{\max}$ ,  $\delta_{i,j} \leq 2\delta_{\max}$ , this leads to  $\delta_{\max} = \epsilon = 1$ . So  $\delta_{i+1,j} = 1$ , which contradicts our hypothesis on  $P_{i+1}$ .
  - (b)  $P_i$  receives a data item, because of Steps 4 through 18, as it is the extremity of a negative maximal slice. Without loss of generality, suppose the slice is  $C_{i,j}$ . By definition of maximal slices, there exists a value  $\epsilon$ , either equal to 0 or 1, such that  $\delta_{i,j} = -2\delta_{\max} + \epsilon$ . Then  $\delta_{i+1,j} = (-2\delta_{\max} + \epsilon) - \delta_{\max} = -3\delta_{\max} + \epsilon$ . As, by definition of  $\delta_{\max}$ ,  $\delta_{i+1,j} \geq -2\delta_{\max}$ , this leads to  $\delta_{\max} = \epsilon = 1$ . So  $\delta_{i,j} = -1$ , which contradicts our hypothesis on  $P_i$ .
  - (c)  $P_i$  receives a data item from another maximal processor, say  $P_{i-1}$ , which sends data items due to Steps 19 and 24. But two maximal processors side by side define a maximal slice. Hence a contradiction because in a maximal slice  $\delta_{i-1,i}$  processor  $P_{i-1}$  sends a data item to processor  $P_{i-2}$  and not to  $P_i$ .
3. After the communications took place, no processor  $P_i$  is such that  $|\delta_i| = \delta_{\max}$ .

As previously, let us consider the case  $\delta_i = \delta_{\max}$  after the communications took place. Because of Case 2, such a case would only arise if the unbalance of  $P_i$  was equal to  $\delta_{\max} - 1$  before the communications (because of the one-port model guaranteed by Lemma 7) and if  $P_i$  received a data item but sent none.

We have three cases to consider:

- (a) Processor  $P_i$  receives a data item from a processor which is the extremity of a maximal slice and which sends data items due to Steps 4 through 18. There is no configuration that can arise where the maximal slice is negative. So the maximal slice is positive. Without any loss of generality, we suppose it is a maximal slice of the form  $C_{i+1,j}$ . Then, by definition of maximal slices,  $\delta_{i+1,j} = 2\delta_{\max} - \epsilon_1$  and  $\epsilon_1$  is either equal to 0 or 1. Thus,  $\delta_{i,j} = 3\delta_{\max} - \epsilon_1 - 1$ . However,  $C_{i,j}$  is not a maximal slice (as, by hypothesis  $P_i$  is not sending any data items). Therefore, by definition of  $\delta_{\max}$ ,  $\delta_{i,j} \leq 2\delta_{\max} - 2$ . Hence, we obtain  $\delta_{\max} \leq \epsilon_1 - 1$  which has no solution.
  - (b) Processor  $P_i$  receives a data item from a processor which is maximal and which sends data items due to Steps 19 through 24. This case can only arise if this maximal processor has a positive unbalance. Without any loss of generality, we suppose that processor  $P_{i-1}$  has an unbalance of  $\delta_{\max}$ . Then,  $\delta_{i-1,i} = 2\delta_{\max} - 1$ . Thus,  $\delta_{i-1,i}$  is a maximal slice, which contradicts the assumption on  $P_{i-1}$ .
  - (c) Processor  $P_i$  receives a data item as it is the extremity of a negative maximal slice, say  $C_{i,j}$ . Then, by definition of maximal slices, there exists  $\epsilon \in \{0, 1\}$  such that  $\delta_{i,j} = -2\delta_{\max} + \epsilon$ . By definition of  $\delta_{\max}$  we have  $\delta_{i+1,j} \geq -2\delta_{\max}$ . As,  $\delta_{i+1,j} = -3\delta_{\max} + \epsilon$ , we obtain  $\delta_{\max} \leq \epsilon$ . Then  $C_{i,j} = -1$  and, because of the test of Step 7, the rules of Steps 8 and 9 do not apply, and  $P_i$  does not receive a data item as it is the extremity of a negative maximal slice.
4. After the communications took place, no non trivial slice  $C_{k,l}$  is such that  $\lceil \frac{|\delta_{k,l}|}{2} \rceil = \delta_{\max}$ .

Once again we only consider the case of positive slices. We can assume that the slice  $C_{k,l}$  was not initially a maximal slice as this case as already been processed. So, there exists a value  $\epsilon_1 \in \{0, 1\}$  such that  $\delta_{k,l} = 2\delta_{\max} - 2 - \epsilon_1$  and we have three cases to consider:

- (a) The slice  $C_{k,l}$  receives a data item from a processor which is the extremity of a maximal slice which sends data items due to Steps 4 through 18. There is no configuration that can arise where the maximal slice is negative. So the maximal slice is positive. Without any loss of generality, we suppose it is of the form  $C_{j,k-1}$ . Then, by definition of maximal slices  $\delta_{j,k-1} = 2\delta_{\max} - \epsilon_2$ , with  $\epsilon_2$  taking values in  $\{0, 1\}$ . Thus,  $\delta_{j,l} = 4\delta_{\max} - \epsilon_1 - \epsilon_2 - 2$ . However, by definition of  $\delta_{\max}$ ,  $\delta_{j,l} \leq 2\delta_{\max}$ . Hence, we obtain  $2\delta_{\max} \leq \epsilon_1 + \epsilon_2 + 2$ . We have two sub-cases to consider:
- i.  $\delta_{\max} = 2$ . Then,  $\epsilon_1 = \epsilon_2 = 1$ . However, in this case  $\delta_{j,l} = 4$ ,  $C_{j,l}$  is a maximal slice, and  $P_l$  sends a data item to  $P_{l+1}$ . Before the communications took place,  $\delta_{k,l} = 1$ . During the communications  $C_{k,l}$  receive at most two data items (as it has two extremities) and send at least one, from  $P_l$ . So, after the communications took place,  $\delta_{k,l}$  is either equal to 0, 1, and 2, and the three cases are fine.
  - ii.  $\delta_{\max} = 1$ . Then, we conclude using the results of Cases 2 and 3.
- (b) The slice  $C_{k,l}$  receives a data item because it is enclosed in a negative maximal slice. Without any loss of generality, we suppose this negative maximal slice is of the form  $C_{k,m}$ .  $C_{k,m}$  being a negative maximal slice,  $\delta_{k,m} = -2\delta_{\max} + \epsilon_2$  with  $\epsilon_2 \in \{0; 1\}$ .
- i.  $l < m$ . Then,  $\delta_{l+1,m} = (-2\delta_{\max} + \epsilon_2) - (2\delta_{\max} - 2 - \epsilon_1) = -4\delta_{\max} + \epsilon_1 + \epsilon_2 + 2$ . By definition of  $\delta_{\max}$ ,  $\delta_{l+1,m} \geq -2\delta_{\max}$ . The case  $\delta_{\max} = 1$  is settled using the result of Case 3. Then  $\delta_{\max} = 2$ ,  $\epsilon_1 = \epsilon_2 = 1$ ,  $\delta_{k,l} = 1$  and  $\delta_{l+1,m} = -4$ . So,  $C_{l+1,m}$  is a negative maximal chain and  $P_l$  sends a data item to  $P_{l+1}$ . So the unbalance of  $C_{k,l}$ , which was originally equal to 1, increases at most by one between before and after communications took place, and there is no problems.
  - ii.  $m < l$ . Then,  $\delta_{m+1,l} = (2\delta_{\max} - 2 - \epsilon_1) - (-2\delta_{\max} + \epsilon_2) = 4\delta_{\max} - \epsilon_1 - \epsilon_2 - 2$ . By definition of  $\delta_{\max}$ ,  $\delta_{m+1,l} \leq 2\delta_{\max}$ . The case  $\delta_{\max} = 1$  is settled using the result of Case 3. Then  $\delta_{\max} = 2$ ,  $\epsilon_1 = \epsilon_2 = 1$ ,  $\delta_{k,l} = 1$  and  $\delta_{m+1,l} = -4$ . So,  $C_{m+1,l}$  is a negative maximal chain and  $P_m$  sends a data item to  $P_{m+1}$ . So the unbalance of  $C_{k,l}$ , which was originally equal to 1, increases at most by one between before and after communications took place, and there is no problems.
- (c) The slice  $C_{k,l}$  only receives a data item from a processor which is maximal and which sends data items because of Steps 19 through 24. This case can only arise if this maximal processor has a positive unbalance. Then, due to Step 19, this is processor  $P_{k-1}$  which has an unbalance of  $\delta_{\max}$ . For  $C_{k,l}$  to be such that  $\lceil \frac{|\delta_{k,l}|}{2} \rceil = \delta_{\max}$  after the communications took place, necessarily,  $\delta_{k,l} \geq 2\delta_{\max} - 2$  before the communications. Then  $\delta_{k-1,l} \geq 3\delta_{\max} - 2$ . As we supposed that  $P_{i-1}$  sends data items because of Steps 19 through 24, the slice  $C_{k-1,l}$  is not maximal and thus  $\delta_{k-1,l} \leq 2\delta_{\max} - 2$ . Hence  $\delta_{\max} \leq 0$ , a contradiction. □

The optimality of Algorithm 3 is a simple corollary of Lemma 8 and of the lower bound defined by Equation 4.

**Theorem 3.** *Algorithm 3 is optimal.*

## 6 Heterogeneous bidirectional ring

In this section, we consider the most general case, that of a heterogeneous bidirectional ring. We do not know any optimal redistribution algorithm in this case. However, if we assume that each processor initially holds more data than it needs to send during the whole execution of algorithm (what we call a *light* redistribution), then we succeed in deriving an optimal solution.

## 6.1 Light redistribution

Throughout this section, we suppose that we have a *light* redistribution: we assume that the number of data items sent by any processor throughout the redistribution algorithm is less than or equal to its original load. There are two reasons for a processor  $P_i$  to send data: (i) because it is overloaded ( $\delta_i > 0$ ); (ii) because it has to forward some data to another processor located further in the ring. If  $P_i$  initially holds at least as many data items as it will send during the whole execution, then  $P_i$  can send at once all these data items. Otherwise, in the general case, some processors may wait to have received data items from a neighbor before being able to forward them to another neighbor.

### 6.1.1 Solution by integer linear programming

Under the “light redistribution” assumption, we can build an integer linear program to solve our problem (see System 6). Let  $\mathcal{S}$  be a solution, and denote by  $\mathcal{S}_{i,i+1}$  the number of data items that processor  $P_i$  sends to processor  $P_{i+1}$ . Similarly,  $\mathcal{S}_{i,i-1}$  is the number of data items that  $P_i$  sends to processor  $P_{i-1}$ . In order to ease the writing of the equations, we impose in the first two equations of System 6 that  $\mathcal{S}_{i,i+1}$  and  $\mathcal{S}_{i,i-1}$  are nonnegative for all  $i$ , which imposes to use other variables  $\mathcal{S}_{i+1,i}$  and  $\mathcal{S}_{i-1,i}$  for the symmetric communications. The third equation states that after the redistribution, there is no more unbalance. We denote by  $\tau$  the execution time of the redistribution. For any processor  $P_i$ , due to the one-port constraints,  $\tau$  must be greater than the time spent by  $P_i$  to send data items (fourth equation) or spent by  $P_i$  to receive data items (fifth equation). Our aim is to minimize  $\tau$ , hence the system:

$$\begin{array}{l} \text{MINIMIZE } \tau, \text{ SUBJECT TO} \\ \left\{ \begin{array}{ll} \mathcal{S}_{i,i+1} \geq 0 & 1 \leq i \leq n \\ \mathcal{S}_{i,i-1} \geq 0 & 1 \leq i \leq n \\ \mathcal{S}_{i,i+1} + \mathcal{S}_{i,i-1} - \mathcal{S}_{i+1,i} - \mathcal{S}_{i-1,i} = \delta_i & 1 \leq i \leq n \\ \mathcal{S}_{i,i+1}c_{i,i+1} + \mathcal{S}_{i,i-1}c_{i,i-1} \leq \tau & 1 \leq i \leq n \\ \mathcal{S}_{i+1,i}c_{i+1,i} + \mathcal{S}_{i-1,i}c_{i-1,i} \leq \tau & 1 \leq i \leq n \end{array} \right. \end{array} \quad (6)$$

**Lemma 9.** *Any optimal solution of System 6 is feasible, for example using the following schedule: for any  $i \in [1, n]$ ,  $P_i$  starts sending data items to  $P_{i+1}$  at time 0 and, after the completion of this communication, starts sending data items to  $P_{i-1}$  as soon as possible under the one-port model.*

*Proof.* We have to show that we are able to schedule the communications defined by any optimal solution  $(\mathcal{S}, \tau)$  of System 6 so that the redistribution takes a time no greater than  $\tau$ . For any  $i \in [1, n]$ , we schedule at time 0, all emissions from  $P_i$  to  $P_{i+1}$ . This communication is done in time  $\mathcal{S}_{i,i+1}c_{i,i+1}$ : because of the “light redistribution” hypothesis,  $P_i$  already holds all the data items that it must send. Because of the fourth equation of System 6, this communication ends before the time  $\tau$ .

For any value of  $i \in [1, n]$ , we still have to schedule the sending of data items from  $P_i$  to  $P_{i-1}$ . We schedule this communication as soon as possible, therefore at time  $\max\{\mathcal{S}_{i,i+1}c_{i,i+1}, \mathcal{S}_{i-2,i-1}c_{i-2,i-1}\}$ , i.e., at the earliest time when (i)  $P_i$  has ended sending data items to  $P_{i+1}$ , and (ii)  $P_{i-1}$  has stopped receiving data items from  $P_{i-2}$ . Therefore, the communication from  $P_i$  to  $P_{i-1}$  ends at the date:

$$\begin{aligned} \max\{\mathcal{S}_{i,i+1}c_{i,i+1}, \mathcal{S}_{i-2,i-1}c_{i-2,i-1}\} + \mathcal{S}_{i,i-1}c_{i,i-1} = \\ \max\{\mathcal{S}_{i,i+1}c_{i,i+1} + \mathcal{S}_{i,i-1}c_{i,i-1}, \mathcal{S}_{i-2,i-1}c_{i-2,i-1} + \mathcal{S}_{i,i-1}c_{i,i-1}\}. \end{aligned} \quad (7)$$

Once again, this is true owing to the “light redistribution” hypothesis: no processor needs to wait to have received some data items before being able to send them to one of its neighbors.

The first term of the “max” expression is the time needed by  $P_i$  to send data items to both  $P_{i+1}$  and  $P_{i-1}$ . This term is less than or equal to  $\tau$  because of the fourth equation of System 6. The second term of the “max” expression is the time needed by  $P_{i-1}$  to receive data items from both  $P_{i-2}$  and  $P_i$ . This term is less than or equal to  $\tau$  because of the fifth equation of System 6.  $\square$

So far, we did not mathematically define a condition for the “light redistribution” hypothesis to hold. In fact, this is not mandatory: we use System 6 to find an optimal solution to the problem. If, in this optimal solution, for any processor  $P_i$ , the total number of data items sent is less than or equal to the initial load ( $\mathcal{S}_{i,i+1} + \mathcal{S}_{i,i-1} \leq L_i$ ), we are under the “light redistribution” hypothesis and we can use the solution of System 6 safely.

### 6.1.2 Solution through rational linear programming

Even if the “light redistribution” hypothesis holds, one may wish to solve the redistribution problem with a technique less expensive than integer linear programming (which is potentially exponential). An idea would be to first solve System 6 to find an optimal *rational* solution, which can always be done in polynomial time, and then to round up the obtained solution to find a “good” integer solution. In fact, the following lemma shows that one of the two natural ways of rounding always lead to an optimal (integer) solution. The complexity of the light redistribution problem is therefore polynomial.

**Proposition 1.** *Let  $\mathcal{R}$  be an optimal rational solution to the redistribution problem. For any  $j$  in  $[1, n]$ ,  $\mathcal{R}_j$  denotes the number of data items that processor  $P_j$  sends to processor  $P_{j+1}$  (using the notations of System 6,  $\mathcal{R}_j = \mathcal{S}_{j,j+1} - \mathcal{S}_{j+1,j}$ ). Let  $\mathcal{F}$  be the integer solution defined by  $\mathcal{F}_1 = \lfloor \mathcal{R}_1 \rfloor$ . Let  $\mathcal{G}$  be the integer solution defined by  $\mathcal{G}_1 = \lceil \mathcal{R}_1 \rceil$ . Then:*

- (i)  $\mathcal{F}$  and  $\mathcal{G}$  are well-defined by the single condition above,
- (ii) either  $\mathcal{F}$  or  $\mathcal{G}$  is an optimal integer solution.

*Proof.* Lemma 10 below states that  $\mathcal{F}$  and  $\mathcal{G}$  are both fully defined. Lemma 11 below states that there exists at least one optimal integer solution  $\mathcal{E}$  such that  $|\mathcal{E}_1 - \mathcal{R}_1| < 1$ . The only two solutions satisfying these constraints are  $\mathcal{F}$  and  $\mathcal{G}$ . Hence the result.  $\square$

**Lemma 10.** *To fully define the number of data items sent between processors in any redistribution scheme, we only need to define, for a single given value of  $j \in [1, n]$ , the number of data items that processor  $P_j$  sends to processor  $P_{j+1}$ .*

*Proof.* Without loss of generality, we suppose we have fixed the value of  $\mathcal{R}_1$ , the number of data items sent by  $P_1$  to  $P_2$ . (Note that  $\mathcal{R}_1$  may be negative, meaning that in fact  $P_2$  sends data items to processor  $P_1$ .) After redistribution, the unbalance of  $P_2$  must be zero. Thus,  $\delta_2 + \mathcal{R}_1 - \mathcal{R}_2 = 0$ . Therefore, as  $\mathcal{R}_1$  is known, the value of  $\mathcal{R}_2$  is also known. Using a direct induction, we then have that, for any value of  $j \in [2, n]$ ,  $\mathcal{R}_j = \delta_j + \mathcal{R}_{j-1}$ , and  $\mathcal{R}_j$  is also known. As  $\sum_{i=1}^n \delta_i = 0$ , one can check that we also have  $\delta_1 + \mathcal{R}_n - \mathcal{R}_1 = 0$ .  $\square$

**Lemma 11.** *Let  $\mathcal{R}$  be an optimal rational solution to the redistribution problem: for any  $j$  in  $[1, n]$ ,  $\mathcal{R}_j$  denotes the number of data items processor  $P_j$  sends to processor  $P_{j+1}$ . Then, there exists an optimal integer solution  $\mathcal{E}$  to the solution problem such that:  $|\mathcal{E}_1 - \mathcal{R}_1| < 1$ .*

*Proof.* We prove Lemma 11 by contradiction. Therefore, we suppose that no optimal integer solution  $\mathcal{E}$  satisfies  $|\mathcal{E}_1 - \mathcal{R}_1| < 1$ . So, let us take an optimal integer solution  $\mathcal{E}$  such that  $|\mathcal{E}_1 - \mathcal{R}_1| \geq 1$ . Let  $\mathcal{R}_1 = \mathcal{E}_1 + z + \epsilon$ , where  $z \in \mathbb{Z}$  and  $\epsilon \in ]-1; 1[$  such that  $\mathcal{E}_1 + z \in [\mathcal{E}_1; \mathcal{R}_1]$ . Therefore

$$\mathcal{E}_1 \leq \mathcal{E}_1 + z \leq \mathcal{E}_1 + z + \epsilon \text{ or } \mathcal{E}_1 \geq \mathcal{E}_1 + z \geq \mathcal{E}_1 + z + \epsilon. \quad (8)$$

Thus, using the construction used in the proof of Lemma 10, we have:

$$\forall i \in [1, n], \mathcal{E}_i \leq \mathcal{E}_i + z \leq \mathcal{E}_i + z + \epsilon \quad \text{or} \quad \forall i \in [1, n], \mathcal{E}_i \geq \mathcal{E}_i + z \geq \mathcal{E}_i + z + \epsilon. \quad (9)$$

Then let  $\mathcal{F}$  be a new integer solution to our problem defined by:  $\mathcal{F}_i = \mathcal{E}_i + z, \forall i \in [1, n]$ . Then,  $|\mathcal{F}_1 - \mathcal{R}_1| = |(\mathcal{E}_1 + z) - (\mathcal{E}_1 + z + \epsilon)| = |\epsilon| < 1$ . If we prove that  $\mathcal{F}$  is an optimal integer solution, we will have reached the desired contradiction.

Consider any value  $i$  in  $[1, n]$ . We have two situations to deal with for processor  $P_i$  (under redistribution  $\mathcal{F}_i$ ):

1.  $\mathcal{F}_{i-1} \cdot \mathcal{F}_i \geq 0$ : under  $\mathcal{F}_i$ , either processor  $P_i$  only communicates data items with one of its neighbors, or it sends data items to one of them and receive data items from the other one.

Without any loss of generality, we suppose that  $\mathcal{F}_{i-1} \geq 0$  and  $\mathcal{F}_i \geq 0$ . Then, we must show that

$$\max\{\mathcal{F}_{i-1}c_{i-1,i}, \mathcal{F}_i c_{i,i+1}\} \leq \tau_{\text{int}},$$

where  $\tau_{\text{int}}$  is the duration of an optimal integer solution. However,  $\mathcal{F}_{i-1}c_{i-1,i} = (\mathcal{E}_{i-1} + z)c_{i-1,i}$ . If  $\mathcal{E}_{i-1} + z$  is null,  $\mathcal{F}_{i-1}c_{i-1,i} = 0 \leq \tau_{\text{int}}$ . Otherwise,  $\mathcal{E}_{i-1} + z$  is not null. As  $\mathcal{E}_{i-1} + z$  is by definition an integer, and as  $|\epsilon| < 1$ ,  $\mathcal{E}_{i-1} + z$  and  $\mathcal{E}_{i-1} + z + \epsilon$  have the same sign, thus are (strictly) positive, and under both redistribution there are data items sent from processor  $P_{i-1}$  to processor  $P_i$ .

- If  $\epsilon > 0$ , then

$$(\mathcal{E}_{i-1} + z)c_{i-1,i} < (\mathcal{E}_{i-1} + z + \epsilon)c_{i-1,i} = \mathcal{R}_{i-1}c_{i-1,i} \leq \tau_{\text{rat}} \leq \tau_{\text{int}},$$

as  $\mathcal{R}$  is by definition an optimal rational solution, and as optimal rational solutions are no worse than optimal integer solutions.

- If  $\epsilon < 0$ , then

$$(\mathcal{E}_{i-1} + z + \epsilon)c_{i-1,i} < (\mathcal{E}_{i-1} + z)c_{i-1,i} < \mathcal{E}_{i-1}c_{i-1,i} \leq \tau_{\text{int}},$$

because of Equation 9, and as  $\mathcal{E}$  is by definition an optimal integer solution.

2.  $\mathcal{F}_{i-1} \cdot \mathcal{F}_i < 0$ : either  $P_i$  receives data items from both of its neighbors, or  $P_i$  sends data items to both of them. Without any loss of generality, we suppose that  $P_i$  sends data items to both of them.

Then, we must show that

$$-\mathcal{F}_{i-1}c_{i,i-1} + \mathcal{F}_i c_{i,i+1} \leq \tau_{\text{int}}. \quad (10)$$

However,  $-\mathcal{F}_{i-1}c_{i,i-1} + \mathcal{F}_i c_{i,i+1} = -(\mathcal{E}_{i-1} + z)c_{i,i-1} + (\mathcal{E}_i + z)c_{i,i+1}$ . As  $\mathcal{E}_{i-1} + z$  is by definition an integer, and as  $|\epsilon| < 1$ ,  $\mathcal{E}_{i-1} + z$  and  $\mathcal{E}_{i-1} + z + \epsilon$  have the same sign, thus are (strictly) negative, and under both redistribution there are data items sent from processor  $P_i$  to processor  $P_{i-1}$ . Similarly, under both redistribution there are data items sent from processor  $P_i$  to processor  $P_{i+1}$ .

As  $\mathcal{R}$  is by definition an optimal rational solution, and as optimal rational solutions are no worse than optimal integer solutions, then:

$$-(\mathcal{E}_{i-1} + z + \epsilon)c_{i,i-1} + (\mathcal{E}_i + z + \epsilon)c_{i,i+1} = -\mathcal{R}_{i-1}c_{i,i-1} + \mathcal{R}_i c_{i,i+1} \leq \tau_{\text{rat}} \leq \tau_{\text{int}}.$$

So, if  $\epsilon(c_{i,i+1} - c_{i,i-1}) \geq 0$ , Equation 10 holds. Otherwise,  $\epsilon(c_{i,i+1} - c_{i,i-1}) < 0$  and we have two cases to consider, depending on the redistribution  $\mathcal{E}$ :

- $\mathcal{E}_{i-1} \cdot \mathcal{E}_i < 0$ : then  $\mathcal{E}_{i-1} < 0$  and  $\mathcal{E}_i > 0$ . Indeed, whatever the redistribution  $\mathcal{S}$  we always have  $\delta_i + \mathcal{S}_{i-1} - \mathcal{S}_i = 0$ . As we have supposed that  $\mathcal{F}_{i-1} < 0$  and  $\mathcal{F}_i > 0$  then  $\delta_i > 0$  which forbids to have  $\mathcal{E}_{i-1} > 0$  and  $\mathcal{E}_i < 0$ .

As  $\mathcal{E}$  is an optimal integer solution, we then have:

$$-\mathcal{E}_{i-1}c_{i,i-1} + \mathcal{E}_i c_{i,i+1} \leq \tau_{\text{int}}.$$

Equation 9 implies that  $z$  and  $\epsilon$  are of same sign. So,  $z(c_{i,i+1} - c_{i,i-1}) < 0$ . Therefore,

$$-\mathcal{F}_{i-1}c_{i,i-1} + \mathcal{F}_i c_{i,i+1} = -(\mathcal{E}_{i-1} + z)c_{i,i-1} + (\mathcal{E}_i + z)c_{i,i+1} < -\mathcal{E}_{i-1}c_{i,i-1} + \mathcal{E}_i c_{i,i+1} \leq \tau_{\text{int}}.$$

- $\mathcal{E}_{i-1} \cdot \mathcal{E}_i \geq 0$ . Without any loss of generality, let us suppose that  $\epsilon > 0$ . Then,  $c_{i,i+1} - c_{i,i-1} < 0$ . Because of Equation 9, as  $\epsilon > 0$  and as  $(\mathcal{E}_{i-1} + z) < 0$ ,  $\mathcal{E}_{i-1} < 0$ , and thus  $\mathcal{E}_i \leq 0$ .

$$\begin{aligned} -(\mathcal{E}_{i-1} + z)c_{i,i-1} + (\mathcal{E}_i + z)c_{i,i+1} &= -\mathcal{E}_{i-1}c_{i,i-1} + \mathcal{E}_i c_{i,i+1} + z(c_{i,i+1} - c_{i,i-1}) \\ &< -\mathcal{E}_{i-1}c_{i,i-1} + \mathcal{E}_i c_{i,i+1}, \end{aligned}$$

as  $c_{i,i+1} - c_{i,i-1} < 0$ . However,  $\mathcal{E}_i \leq 0$ , so

$$-(\mathcal{E}_{i-1} + z)c_{i,i-1} + (\mathcal{E}_i + z)c_{i,i+1} < -\mathcal{E}_{i-1}c_{i,i-1} \leq \tau_{\text{int}}$$

as  $\mathcal{E}$  is an optimal integer solution. □

## 6.2 General case

### 6.2.1 Lower bound

We have the following bound on the optimal redistribution time:

**Lemma 12.** *Let  $\tau$  be the optimal redistribution time. Then:*

$$\tau \geq \max \left\{ \begin{array}{l} \max_{1 \leq k \leq n, \delta_k > 0} \delta_k \min\{c_{k,k-1}, c_{k,k+1}\}, \\ \max_{1 \leq k \leq n, \delta_k < 0} -\delta_k \min\{c_{k-1,k}, c_{k+1,k}\}, \\ \max_{\substack{1 \leq k \leq n, \\ 1 \leq l \leq n-2, \\ \delta_{k,k+l} > 0}} \min_{0 \leq i \leq \delta_{k,k+l}} \max\{i \cdot c_{k,k-1}, (\delta_{k,k+l} - i) \cdot c_{k+l,k+l+1}\} \\ \max_{\substack{1 \leq k \leq n, \\ 1 \leq l \leq n-2, \\ \delta_{k,k+l} < 0}} \min_{0 \leq i \leq -\delta_{k,k+l}} \max\{i \cdot c_{k-1,k}, (-\delta_{k,k+l} - i) \cdot c_{k+l+1,k+l}\} \end{array} \right\} \quad (11)$$

*Proof.* Consider any processor  $P_i$  with positive unbalance ( $\delta_i > 0$ ). Even if processor  $P_i$  can send data items to both of its neighbors, because of the one-port model, it cannot send data items to both of them *simultaneously*. The best way for processor  $P_i$  to send  $\delta_i$  data items is then to send them using the fastest of its outgoing links. So, it requires processor  $P_i$  at least a time of  $\delta_i \times \min\{c_{i,i-1}, c_{i,i+1}\}$  to send  $\delta_i$  data items, whatever the destinations of these data items. We have a symmetric result for the case  $\delta_i < 0$ . Hence the first two equations of the System 11.

Now, consider any non trivial slice of consecutive processors  $C_{k,l}$ . By “non trivial” we mean that the slice is not reduced to a single processor (we already treated that case) and that it does not contain all processors. We suppose that  $\delta_{k,l} > 0$ . So, in any redistribution scheme, at least  $\delta_{k,l}$  data items must be sent by  $C_{k,l}$ . As this slice is not reduced to a single processor, the two processors at the extremities of the slice,  $P_k$  and  $P_l$ , can simultaneously send data items to their neighbors outside of the slice,  $P_{k-1}$  and  $P_{l+1}$  respectively. Therefore, during the redistribution, processor  $P_k$  sends a certain amount  $i \in [0, \delta_{k,l}]$  of data items to processor  $P_{k-1}$ , while processor  $P_l$  sends the remaining data items to  $P_{l+1}$ , which takes a time  $\max\{i \cdot c_{k,k-1}, (\delta_{k,l} - i) \cdot c_{l,l+1}\}$ . Then we chose for  $i$  a value which minimizes this time. We have a symmetric result for the case  $\delta_{k,l} < 0$ . Hence the last two equations of the System 11. □

### 6.2.2 Heuristic approaches

We do not know whether the bound given by Lemma 12 can always be reached, but we have no counter-example proving that the bound is not tight.

When the solution found by System 6 does not satisfy the “light redistribution” hypothesis, there is the possibility to modify the system to enforce it: we obtain System 12 which finds a

solution which satisfies the “light redistribution” hypothesis, if one exists. But there is no reason *a priori* for the solution of System 12 to be optimal.

$$\begin{array}{l} \text{MINIMIZE } \tau, \text{ SUBJECT TO} \\ \left\{ \begin{array}{ll} \mathcal{S}_{i,i+1} \geq 0 & 1 \leq i \leq n \\ \mathcal{S}_{i,i-1} \geq 0 & 1 \leq i \leq n \\ \mathcal{S}_{i,i+1} + \mathcal{S}_{i,i-1} - \mathcal{S}_{i+1,i} - \mathcal{S}_{i-1,i} = \delta_i & 1 \leq i \leq n \\ \mathcal{S}_{i,i+1}c_{i,i+1} + \mathcal{S}_{i,i-1}c_{i,i-1} \leq \tau & 1 \leq i \leq n \\ \mathcal{S}_{i+1,i}c_{i+1,i} + \mathcal{S}_{i-1,i}c_{i-1,i} \leq \tau & 1 \leq i \leq n \\ \mathcal{S}_{i,i+1} + \mathcal{S}_{i,i-1} \leq L_i & 1 \leq i \leq n \end{array} \right. \end{array} \quad (12)$$

To conclude this section, we point out that the design of an optimal algorithm in the most general case remains open. Given the complexity of the lower bound, the problem looks very difficult to solve.

## 7 Related work

Redistribution algorithms have been the focus of an abundant literature. On the theoretical side, in the framework of High Performance Fortran [25] compilation, Kremer [26] showed the NP-completeness of a simple redistribution problem. This negative results shows that optimal algorithms can be designed only for particular cases, such as the ring architecture in this paper. To the best of our knowledge, no other redistribution algorithms has been proven optimal, but several efficient algorithms have been designed for rings [20, 28, 13], trees or hypercubes [41]. The elastic load balancing algorithm designed in [30, 4] has led to a data redistribution software used for query processing [8] and medical image analysis [35].

The block-cyclic distribution of data arrays plays a very important role in scientific libraries [5]. In a **CYCLIC**( $r$ ) distribution over  $p$  processors, blocks of  $r$  consecutive elements of the array are distributed to the processors in a wraparound fashion, and the parameter  $r$  is chosen to optimize the granularity, i.e. the computation-to-communication ratio. Because this granularity changes from one computational kernel to the other, moving from a **CYCLIC**( $r$ ) distribution over  $p$  processors to a **CYCLIC**( $s$ ) distribution over  $q$  processors is a very useful redistribution procedure, which has been implemented using a caterpillar algorithm in ScaLAPACK [34]. Several papers, including [23, 39, 14, 33, 19, 11, 24], have dealt with various optimizations of this redistribution procedure. Along this line of research, automatic data redistribution tools are presented in [19].

Even though we did not deal with load-balancing algorithms in this paper, we quote some key references on the subject. For homogeneous platforms, see the collection of papers [38], and for heterogeneous clusters see chapter 25 in [9]. Several authors [17, 32, 31, 40, 21] propose a mapping policy which dynamically minimizes system degradation (including the cost of remapping) for each computation step. Static strategies aiming at distributing independent chunks of work to two-dimensional processor grids are studied in [1, 2]. Relaxing the geometrical constraints induced by two-dimensional grids leads to irregular partitionings [12, 22, 3] that allow for a good load-balancing but are much more difficult to implement. This approach has been extended to three-dimensional problems [18].

Finally, we briefly mention three sample applications whose implementation can directly benefit from the redistribution strategies designed in this paper. The analysis of pulses propagating in a nonlinear medium calls for adaptive computational windows, and redistribution must occur frequently as the computation progresses [6]. A two-level redistribution procedure is advocated in [27] for structured adaptive mesh refinement. A multi-level diffusion re-partitioner is presented in [36, 37] for irregular grid computations and has been incorporated into the PARMETIS library. Of course this short list could be extended dramatically.



## 8 Experimental results

To evaluate the impact of the redistributions, we used the SIMGRID [29] simulator to model an iterative application, implemented on a platform generated with the Tiers network generator [10, 15].

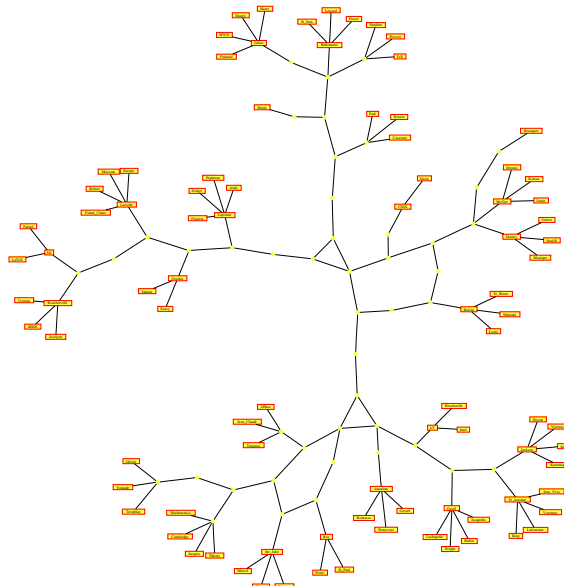


Figure 11: The platform is composed of 90 machine nodes, connected through 192 communication links.

We use the platform represented in Figure 11. The capacities of the edges are assigned using the classification of the Tiers generator (local LAN link, LAN/MAN link, MAN/WAN link, . . .). For each link type, we use values measured using `pathchar` [16] between some machines in ENS Lyon and some other machines scattered in France (Strasbourg, Lille, Grenoble, and Orsay), in the USA (Knoxville, San Diego, and Argonne), and in Japan (Nagoya and Tokyo).

We randomly select  $p$  processors in the platform to build the execution ring. The communication speed is given by the slowest link in the route from a processor to its successor (or predecessor) in the ring. The processing powers (CPU speeds) of the nodes are first randomly chosen in a list of values corresponding to the processing powers (expressed in MFlops and evaluated thanks to a benchmark taken from LINPACK [7]) of a wide variety of machines (Pentium Pro 200MHz, Pentium 2 350MHz, Celeron 400MHz, Athlon 1.4GHz, Pentium 4 1.7GHz, . . .). But we make these speeds vary during the execution of the application.

We model an iterative application which executes during 100 iterations. At each iteration, independent data are updated by the processors. We may think of a  $m \times n$  data matrix whose columns are distributed to the processors (we use  $n = m = 1000$  in the experiment). Ideally, each processor should be allocated a number of columns proportional to its CPU speed. This is how the distribution of columns to processors is initialized.

To motivate the need for redistributions, we create an unbalance by letting the CPU speeds vary during the execution. The speed of each processor changes two times, first at some iteration randomly chosen between iterations number 20 and 40, and then at some iteration randomly chosen between iterations number 60 and 80) for each node to change the processing power (see Figure 12 for an illustration). We record the values of each CPU speed in a SIMGRID trace.

In the simulations, we use the heterogeneous bidirectional algorithm for light redistributions, and we test five different schemes, each with a given number of redistributions within the 100 iterations. The first scheme has no redistribution at all. The second scheme implements a redistribution after iteration number 50. The third scheme uses four redistributions, after iterations

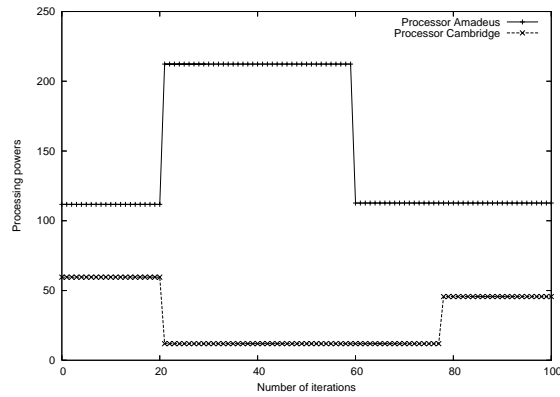


Figure 12: Processing power of 2 sample machine nodes.

20, 40, 60 and 80. The fourth scheme uses 9 redistributions, implemented every 10 iterations, and the last one uses 19 redistributions, implemented every 5 iterations. Given the shape of the CPU traces, some redistributions are likely to be beneficial during the execution.

The last parameter to set is the computation-to-communication ratio, which amounts to set the relative (average) cost of a redistribution versus the cost of an iteration. When this parameter increases, iterations take more time, and the usefulness of a redistribution becomes more important.

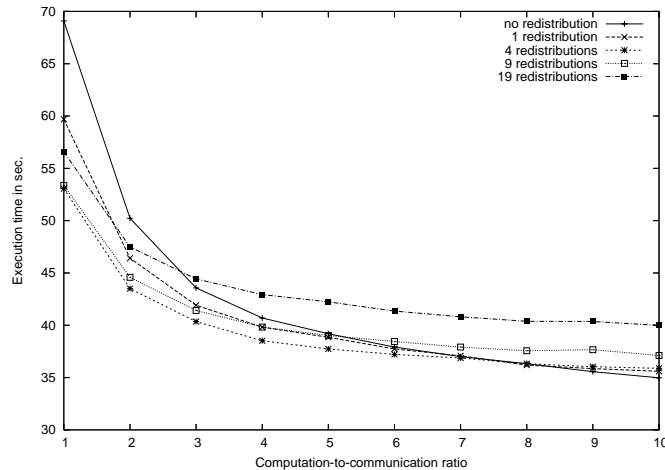


Figure 13: Normalized execution time as a function of the computation-to-communication ratio, for a ring of 8 processors.

In Figures 13 and 14, we plot the execution time of different computation schemes. Both figures report the same comparisons, but for different ring sizes: we use 8 processors in Figures 13, and 32 in Figures 14.

As expected, when the processing power is high (ratio = 10 in the figures), the best strategy is to use no redistribution, as their cost is prohibitive. Conversely, when the processing power is low (ratio = 1 in the figures), it pays off to use many redistributions, but not too many! As the ratio increases, all tradeoffs can be found.

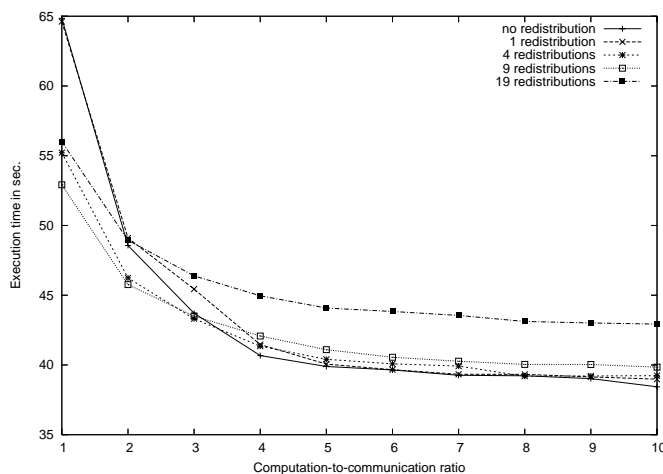


Figure 14: Normalized execution time as a function of the ratio computation-to-communication, for a ring of 32 processors.

## 9 Conclusion

In this paper, we have considered the problem of redistributing data on rings of processors. For homogeneous rings the problem has been completely solved. Indeed, we have designed optimal algorithms, and provided formal proofs of correctness, both for unidirectional and bidirectional rings. The bidirectional algorithm turned out to be quite complex, and requires a lengthy proof.

For heterogeneous rings there remains further research to be conducted. The unidirectional case was easily solved, but the bidirectional case remains open. Still, we have derived an optimal solution for light redistributions, an important case in practice. The complexity of the bound provided for the general case shows that designing an optimal algorithm is likely to be a difficult task.

All our algorithms have been implemented and extensively tested. We have reported some simulation results for the most difficult combination, that of heterogeneous bi-directional rings. As expected, the cost of data redistributions may not pay off a little unbalance of the work in some cases. Further work will aim at investigating how frequently redistributions must occur in real-life applications.

## References

- [1] J. Barbosa, J. Tavares, and A. J. Padilha. Linear algebra algorithms in a heterogeneous cluster of personal computers. In *9th Heterogeneous Computing Workshop (HCW'2000)*, pages 147–159. IEEE Computer Society Press, 2000.
- [2] O. Beaumont, V. Boudet, A. Petitet, F. Rastello, and Y. Robert. A proposal for a heterogeneous cluster ScaLAPACK (dense linear solvers). *IEEE Trans. Computers*, 50(10):1052–1070, 2001.
- [3] O. Beaumont, V. Boudet, F. Rastello, and Y. Robert. Matrix multiplication on heterogeneous platforms. *IEEE Trans. Parallel Distributed Systems*, 12(10):1033–1051, 2001.
- [4] A. Bevilacqua. A dynamic load balancing method on a heterogeneous cluster of workstations. *Informatica*, 23(1):49–56, 1999.
- [5] L. S. Blackford, J. Choi, A. Cleary, E. D’Azevedo, J. Demmel, I. Dhillon, J. Dongarra, S. Hammarling, G. Henry, A. Petitet, K. Stanley, D. Walker, and R. C. Whaley. *ScaLAPACK Users’ Guide*. SIAM, 1997.

- [6] A. Bourgade and B. Nkonga. Dynamic load balancing computation of pulses propagating in a nonlinear medium. *The Journal of Supercomputing*, 28(3):279–294, 2004.
- [7] R. P. Brent. The LINPACK Benchmark on the AP1000: Preliminary Report. In *CAP Workshop 91*. Australian National University, 1991. Website <http://www.netlib.org/linpack/>.
- [8] L. Brunie, A. Flory, and H. Kosch. New static scheduling and elastic load balancing methods for parallel query processing. In *Basque International Workshop on Information Technology BIWIT*. IEEE Computer Society Press, 1995.
- [9] R. Buyya. *High Performance Cluster Computing. Volume 1: Architecture and Systems*. Prentice Hall PTR, Upper Saddle River, NJ, 1999.
- [10] K. L. Calvert, M. B. Doar, and E. W. Zegura. Modeling internet topology. *IEEE Communications Magazine*, 35(6):160–163, June 1997. Available at <http://citeseer.nj.nec.com/calvert97modeling.html>.
- [11] C.H.Hsu, Y. Chung, D. Yang, and C. Dow. A generalized processor mapping technique for array redistribution. *IEEE Trans. Parallel Distributed Systems*, 12(7):743–757, 2001.
- [12] P. E. Crandall and M. J. Quinn. Block data decomposition for data-parallel programming on a heterogeneous workstation network. In *2nd International Symposium on High Performance Distributed Computing*, pages 42–49. IEEE Computer Society Press, 1993.
- [13] E. Deelman and B. Szymanski. Dynamic load balancing in parallel discrete event simulation for spatially explicit problems. In *PADS'98, 12th Workshop on Parallel and Distributed Simulation*, pages 46–53. IEEE Computer Society Press, 1998.
- [14] F. Desprez, J. Dongarra, A. Petitet, C. Randriamaro, and Y. Robert. Scheduling block-cyclic array redistribution. *IEEE Trans. Parallel Distributed Systems*, 9(2):192–205, 1998.
- [15] M. Doar. A better model for generating test networks. In *Proceedings of Globecom '96*, Nov. 1996. Available at <http://citeseer.nj.nec.com/doar96better.html>.
- [16] A. B. Downey. Using pathchar to estimate internet link characteristics. In *Measurement and Modeling of Computer Systems*, pages 222–223, 1999. Available at <http://citeseer.nj.nec.com/downey99using.html>.
- [17] J. E. Flaherty, R. M. Loy, C. Özturan, M. S. Shephard, B. K. Szymanski, J. D. Teresco, and L. H. Ziantz. Parallel structures and dynamic load balancing for adaptive finite element computation. *Applied Numerical Mathematics*, 26(1-2):241–263, 1997.
- [18] J. E. Flaherty, R. M. Loy, M. S. Shephard, B. K. Szymanski, J. D. Teresco, and L. H. Ziantz. Adaptive local refinement with octree load balancing for the parallel solution of three-dimensional conservation laws. *J. Parallel and Distributed Computing*, 47(2):139–152, 1997.
- [19] J. Garcia, E. Ayguadé, and J. Labarta. A framework for integrating data alignment, distribution, and redistribution in distributed memory multiprocessors. *IEEE Trans. Parallel Distributed Systems*, 12(4):416–431, 2001.
- [20] M. Hamdi and C. Lee. Dynamic load balancing of data parallel applications on a distributed network. In *9th International Conference on Supercomputing ICS'95*, pages 170–179. ACM Press, 1995.
- [21] Y. Hu and R. Blake. Load balancing for unstructured mesh applications. *Parallel and Distributed Computing Practices*, 2(3), 1999.
- [22] M. Kaddoura, S. Ranka, and A. Wang. Array decomposition for nonuniform computational environments. *Journal of Parallel and Distributed Computing*, 36:91–105, 1996.

- [23] E. T. Kalns and L. M. Ni. Processor mapping techniques towards efficient data redistribution. *IEEE Trans. Parallel Distributed Systems*, 6(12):1234–1247, 1995.
- [24] J. Knoop and E. Mehofer. Distribution assignment placement: effective optimization of redistribution costs. *IEEE Trans. Parallel Distributed Systems*, 13(6):628–647, 2002.
- [25] C. H. Koelbel, D. B. Loveman, R. S. Schreiber, G. L. S. Jr., and M. E. Zosel. *The High Performance Fortran Handbook*. The MIT Press, 1994.
- [26] U. Kremer. NP-Completeness of dynamic remapping. In *Proceedings of the Fourth Workshop on Compilers for Parallel Computers*, Delft, The Netherlands, 1993. also available as Rice Technical Report CRPC-TR93330-S.
- [27] Z. Lan, V. Taylor, and G. Bryan. Dynamic load balancing of samr applications on distributed systems. In *Proceedings of the ACM/IEEE Symposium on Supercomputing (SC'01)*. IEEE Computer Society Press, 2001.
- [28] C. Lee and M. Hamdi. Parallel image processing applications on a network of workstations. *Parallel Computing*, 21:137–160, 1995.
- [29] A. Legrand, L. Marchal, and H. Casanova. Scheduling Distributed Applications: The SIM-GRID Simulation Framework. In *Proceedings of the Third IEEE International Symposium on Cluster Computing and the Grid (CCGrid'03)*, May 2003.
- [30] S. Miguet and Y. Robert. Elastic load balancing for image processing algorithms. In H. Zima, editor, *Parallel Computation*, LNCS 591, pages 438–451. Springer Verlag, 1992.
- [31] D. Nicol and J. P.F. Reynolds. Optimal dynamic remapping of data parallel computations. *IEEE Trans. Computers*, 39(2):206–219, 1990.
- [32] D. Nicol and J. Saltz. Dynamic remapping of parallel computations with varying resource demands. *IEEE Trans. Computers*, 37(9):1073–1087, 1988.
- [33] N. Park, V. Prasanna, and C. Raghavendra. A framework for integrating data alignment, distribution, and redistribution in distributed memory multiprocessors. *IEEE Trans. Parallel Distributed Systems*, 10(12):1217–1240, 1999.
- [34] L. Prylli and B. Tourancheau. Fast runtime block-cyclic data redistribution on multiprocessors. *J. Parallel Distributed Computing*, 45:63–72, 1997.
- [35] D. Sarrut and S. Miguet. ARAMIS: a remote access medical imaging system. In *ISCOPE'99, 3rd International Symposium on Computing in Object-Oriented Parallel Environments*, volume 1732 of *Lecture Notes in Computer Science*. Springer, 1999.
- [36] K. Schloegel, G. Karypis, and V. Kumar. Multilevel diffusion schemes for repartitioning of adaptive meshes. volume 47, pages 109–124, 1997.
- [37] K. Schloegel, G. Karypis, and V. Kumar. A unified algorithm for load-balancing adaptive scientific simulations. In *Proceedings of the ACM/IEEE Symposium on Supercomputing (SC'00)*. IEEE Computer Society Press, 2000.
- [38] B. A. Shirazi, A. R. Hurson, and K. M. Kavi. *Scheduling and load balancing in parallel and distributed systems*. IEEE Computer Science Press, 1995.
- [39] R. Thakur, A. Choudhary, and J. Ramanujam. Efficient algorithms for array redistribution. *IEEE Trans. Parallel and Distributed Systems*, 7(6):587–594, 1996.
- [40] J. Watts and S. Taylor. A practical approach to dynamic load balancing. *IEEE Trans. Parallel and Distributed Systems*, 9(93):235–248, 1998.
- [41] M.-Y. Wu. On runtime parallel scheduling for processor load balancing. *IEEE Trans. Parallel and Distributed Systems*, 8(2):173–186, 1997.