

# Dynamic Load-Balancing with Variable Number of Processors based on Graph Repartitioning

Clément Vuchener, Aurélien Esnard

► **To cite this version:**

Clément Vuchener, Aurélien Esnard. Dynamic Load-Balancing with Variable Number of Processors based on Graph Repartitioning. [Research Report] RR-7926, INRIA. 2012. <hal-00687073>

**HAL Id: hal-00687073**

**<https://hal.inria.fr/hal-00687073>**

Submitted on 12 Apr 2012

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

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



# Dynamic Load-Balancing with Variable Number of Processors based on Graph Repartitioning

Clément Vuchener, Aurélien Esnard

**RESEARCH  
REPORT**

**N° 7926**

April 2012

Project-Team HiePACS





## Dynamic Load-Balancing with Variable Number of Processors based on Graph Repartitioning

Clément Vuchener\*, Aurélien Esnard\*

Project-Team HiePACS

Research Report n° 7926 — April 2012 — 14 pages

**Abstract:** Dynamic load balancing is an important step conditioning the performance of parallel adaptive codes whose load evolution is difficult to predict. Most of the works which answer this problem perform well, but are limited to an initially fixed number of processors which is not modified at runtime. These approaches can be very inefficient, especially in terms of resource consumption. In this paper, we present a new graph repartitioning algorithm which accepts a variable number of processors, assuming the load is already balanced. Our algorithm minimizes both data communication and data migration overheads, while maintaining the computational load balanced. This algorithm is based on a theoretical result, that constructs optimal communication matrices with both a minimum migration volume and a minimum number of communications. An experimental study which compares our work against state-of-the-art approaches is presented.

**Key-words:** dynamic load-balancing, graph, hypergraph, partitioning, migration, repartitioning

\* Univ. Bordeaux, LaBRI, UMR 5800, HiePACS Project, INRIA, F-33400 Talence, France.

**RESEARCH CENTRE  
BORDEAUX – SUD-OUEST**

351, Cours de la Libération  
Bâtiment A 29  
33405 Talence Cedex

## Équilibrage de Charge Dynamique avec Nombre Variable de Processeurs basé sur un Repartitionnement de Graphe

**Résumé :** L'équilibrage de charge dynamique est une étape importante qui conditionne la performance des codes parallèles adaptatifs, dont l'évolution de la charge est difficile à prédire. La plupart des travaux qui s'intéressent à ce problème sont aujourd'hui matures, mais se limitent au cas où le nombre de processeurs initialement fixé ne change pas lors de l'exécution. Ces approches peuvent être particulièrement inefficaces en termes de consommation des ressources. Dans ce papier, nous présentons un nouvel algorithme de repartitionnement de graphe, qui accepte un nombre variable de processeurs, en supposant que la charge est déjà équilibrée. Notre algorithme minimise conjointement les temps de communication et de migration des données, tout en maintenant la charge équilibrée. Cet algorithme s'appuie sur un résultat théorique, qui permet de construire des matrices de communication optimales, minimisant à la fois le volume de migration et le nombre de messages échangés. Nous validons notre approche par une étude expérimentale comparative avec les méthodes classiques.

**Mots-clés :** équilibrage de charge dynamique, graphe, hypergraphe, partitionnement, migration, repartitionnement

## 1 Introduction

In the field of scientific computing, the load-balancing is a crucial issue, which determines the performance of parallel programs. As a general rule, one applies a static balancing algorithm, which equilibrates the computational load between processors before running the parallel program. For some scientific applications, such as adaptive codes (e.g. adaptive mesh refinement), the evolution of the load is unpredictable. Therefore, it is required to periodically compute a new balancing at runtime, using a dynamic load-balancing algorithm. As this step may be performed frequently, it must use a fast and incremental algorithm with a quality trade-off. As computation progresses, the global workload may increase drastically, exceeding memory limit for instance. In such a case, we argue it should be relevant to adjust the number of processors while maintaining the load balanced. However, this is still an open question that we investigate in this paper.

A very common approach to solve the load-balancing problem (static or dynamic) is based on graph model. Each vertex of the graph represents a basic computational task and each edge represents a dependency in the calculation between two tasks. To equilibrate the charge between  $M$  processors, one performs a *graph partitioning* in  $M$  parts, each part being assigned to a given processor. More precisely, the objective consists of dividing the graph into  $M$  parts (or vertex subsets), such that the parts are disjoint and have equal size, and there are few edges cut between the parts. Here are the classical partitioning criteria:

- minimize the computation time ( $T_{comp}$ ), which consists of dividing the graph in parts of equal weight (up to an unbalance factor);
- minimize the communication time ( $T_{comm}$ ), which consists of minimizing the edge cut of the graph induced by the new partition.

If the load changes at runtime, the current partition becomes unbalanced and it is required to perform a *graph repartitioning*. In addition to the classical partitioning criteria, the problem of repartitioning optimizes the following criteria [2]:

- minimize the migration time ( $T_{mig}$ ), which consists of minimizing the vertex weight moving from the old partition to the new one;
- minimize the repartitioning time ( $T_{repart}$ ).

It should be noticed that the repartitioning and migration steps are not performed at each iteration in the application, but periodically (e.g. every  $\alpha$  iterations). As a consequence, the total time period of the code is written:  $T_{total} = \alpha.(T_{comp} + T_{comm}) + T_{mig} + T_{repart}$ . Assuming  $T_{repart}$  is negligible compared to the other terms, and if we consider that  $T_{comp}$  is implicitly minimized by balancing the parts, it follows that to minimize  $T_{total}$ , one must minimize  $\alpha.T_{comm} + T_{mig}$ . Finally, it clearly shows there is a trade-off between the optimization of the communication time ( $T_{comm}$ ) and optimization of the migration time ( $T_{mig}$ ). This compromise is controlled by the parameter  $\alpha$ , which depends on the target application.

As we will see in the following section, there are many works around the dynamic load-balancing and graph repartitioning. However, all these works are limited—as far as we know—to the case where the number of processors is initially fixed and will not be modified at runtime. This can be very inefficient, especially in terms of resource consumption [8]. To overcome this issue, we propose at section 3 a new graph repartitioning algorithm which accepts a variable number of processors, assuming the load is already balanced. We call this problem the  $M \times N$  graph repartitioning problem. Our algorithm minimizes both data communication (i.e. cut size) and data migration overheads, while maintaining the computational load balance in parallel. This algorithm is based on a theoretical result, that constructs optimal communication matrices with

both a minimum migration volume and a minimum number of communications (see Sec. 3.2). Moreover, it uses recent graph partitioning technique with fixed vertices to take into account migration constraints. Finally, we validate this algorithm at section 4 with some experimental results, that compare our approach with state-of-the-art partitioning softwares.

## 2 Related Work

There are many works in the field of dynamic load-balancing [6, 13]. We briefly review the most popular methods based on graph (or hypergraph) repartitioning techniques.

The simplest approach is certainly the *Scratch-Remap* scheme [10], which calculates a new partitioning from scratch, that is to say, without taking into account the old partition. This technique obviously minimizes the cut, but does not control the migration at all. To reduce this latter cost, an additional step of *remapping* attempts to renumber the new parts in order to maximize data remaining in place.

Another approach are the diffusive methods. In their simplest form, they are based on the heat equation to dynamically equilibrate the load [12]. It is an interactive algorithm, where two neighboring processors exchange at each step an amount of data proportional to their load difference. After several steps, the convergence of the diffusion scheme reaches a new load balancing, that defines a new partitioning.

A more recent approach consists in repartitioning graph (or hypergraph) by minimizing both the cut size and the data movement due to migration (*RM-Metis* [1] and *Zoltan* [2]). For each part, a *fixed vertex* of zero weight is added. This particular vertex is connected by new edges—called *migration edges*—to all regular vertices that corresponds to this part. Then, one performs a partitioning of this enriched graph, with the constraint that fixed vertices are required to be assigned to their respective part in the final solution. Other vertices are free to move. Thus, if a normal vertex changes its part, this involves to cut a migration edge and to pay for an additional migration cost associated with this edge. As a partitioner attempt to minimize the cut size, it will also minimize the data movement due to migration. *Scotch* has recently added a similar graph repartitioning method based on fixed vertices, using a local diffusive refinement [4].

One can find in the literature many other works on dynamic load-balancing, including geometric methods like *Recursive Coordinate Bisection* (RCB) [6] or *Space-Filling Curve* (SFC) [11], spectral methods [14], or still more exotic approaches such as *skewed graph* partitioning [7].

All these works are very interesting, but are limited to the case where the number of processors is initially fixed and is not modified at runtime. In our knowledge, there is no research that investigates the problem of graph (or hypergraph) partitioning with a variable number of processors. However, some recent studies have shown the interest of such an approach, by dynamically adjusting the number of processors in an adaptive code (AMR) to optimize both the parallel runtime and resource consumption [8].

## 3 Algorithm

This section presents our graph repartitioning algorithm which accepts a variable number of processors, assuming the load is already balanced. It is based on a theoretical result on *optimal* communication matrices, that minimizes both the data volume and the number of communications during the migration phase. These matrices are conveniently represented by a *repartitioning hypergraph*, that captures the optimal communication scheme we will impose. Then, the initial graph is enriched with fixed vertices, that models our additional migration constraints in a similar way to *Zoltan* [2] or *RM-Metis* [1]. Thus, the partitioning of this graph will minimize both

the regular cut size and the data movement due to migration, while respecting the optimal communication scheme.

### 3.1 Communication Matrix and Repartitioning Hypergraph

Let consider a graph  $G = (V, E)$ , where  $V$  is the set of vertices, and  $E$  is the set of edges. Let  $w$  be the weight function that maps to a vertex subset of  $G$  its weight. We notice  $W = w(V)$  the weight of the whole graph. Let  $P = (V_1, V_2, \dots, V_M)$  be the initial partition of  $V$  into  $M$  parts and  $P' = (V'_1, V'_2, \dots, V'_N)$  the final partition into  $N$  parts.

Let  $C = (C_{i,j})$  be the  $M \times N$  communication matrix associated with the repartitioning of  $G$  from  $P$  to  $P'$ . The element  $C_{i,j}$  is the amount of data sent by the processor  $i$  to the processor  $j$ . According to the graph model,  $C_{i,j}$  is equal to  $w(V_i \cap V'_j)$ . In this paper, we focus on *perfect* communication matrix, which results from two perfectly balanced partitions,  $P$  and  $P'$ . Such matrices satisfy the following constraints: for each row  $i$ ,  $w(V_i) = \sum_{1 \leq j \leq N} C_{i,j} = W/M$  (row constraint) and for each column  $j$ ,  $w(V'_j) = \sum_{1 \leq i \leq M} C_{i,j} = W/N$  (column constraint). As a consequence,  $W$  must be a multiple of both  $M$  and  $N$ .

We define the *number of communications*,  $Z(C)$ , as the number of non-zero terms in  $C$ . It represents the number of messages exchanged between former and newer parts, including “in-place” communications from a processor to itself. In the case of perfect communication matrix, we will demonstrate in the following section that this number is minimum for  $M + N - GCD(M, N)$  and obviously maximum for  $M.N$ . Then, we define the *migration volume*,  $Mig(C)$ , as the amount of data being sent to a different processor, *i.e.*  $Mig(C) = \sum_{i \neq j} C_{i,j}$ .

The matrix  $C$  can be interpreted as an hypergraph  $H$ , called *repartitioning hypergraph*. This hypergraph is composed of  $M$  vertices representing the initial parts and  $N$  hyperedges representing the new parts obtained after the repartitioning step. A vertex  $i$  of  $H$  belongs to an hyperedge  $j$  if data are exchanged between the old part  $i$  and the new part  $j$  during the migration. The repartitioning hypergraph allows to model the *communication scheme* without detailing the volume of data exchanged as the communication matrix does. We will see how this hypergraph representation makes easier to solve the correspondence problem we have in section 3.3.

### 3.2 Optimal Communication Matrices

Our goal in this section is to seek communication matrices with good properties to perform efficiently the migration step. To simplify our discussion, we will assume in all this section that the communication matrix  $C$  of dimension  $M \times N$  is perfect with  $W = M.N$ . As the initial and final partition are perfectly balanced, a source processor sends a data volume of  $N$  and a target processor receives a data volume of  $M$  (including “in-place” communications).

**Definition 1.** *In this paper, a perfect communication matrix  $C$  is said to be optimal if it minimizes both the migration volume  $Mig(C)$  and the number of communications  $Z(C)$ .*

**Theorem 1.** *Let  $C$  be a perfect communication matrix of dimension  $M \times N$ . The minimum number of communications is  $M + N - GCD(M, N)$ .*

*Proof.* Let  $\mathcal{G} = ((A, B), E)$  be the bipartite graph that represents the communication of matrix  $C$  from  $M = |A|$  processors to  $N = |B|$  processors. Let  $K$  be the number of connected components of  $\mathcal{G}$ , noted  $\mathcal{G}_i = ((A_i, B_i), E_i)$  with  $1 \leq i \leq K$ . For each component  $\mathcal{G}_i$ ,  $M_i = |A_i|$  processors send a data volume  $M_i.N$  to  $N_i = |B_i|$  processors that receive a data volume  $N_i.M$ . Therefore,  $\mathcal{G}_i$  exchange a data volume  $V_i = M_i.N = N_i.M$ , with  $M_i$  and  $N_i$  non null. As  $V_i$  is multiple of both  $M$  and  $N$ , one can say  $V_i \geq LCM(M, N)$ . Consequently, the total volume of communications  $M.N = \sum_{i \in [1, K]} V_i$  is superior or equal to  $K.LCM(M, N)$ . As  $GCD(M, N).LCM(M, N) =$



$M.N$ , one can deduce  $K \leq GCD(M, N)$ . As  $\mathcal{G}_i$  is a connected graph, its number of edges  $|E_i|$  is superior or equal to  $M_i + N_i - 1$ . And the total number of edges  $|E| = \sum_{i \in [1, k]} |E_i|$  is superior or equal to  $\sum_{i \in [1, K]} M_i + \sum_{i \in [1, K]} N_i - K = M + N - K$ . As a consequence, the total number of communications  $|E|$  is superior or equal to  $M + N - GCD(M, N)$ , for  $K \leq GCD(M, N)$ .  $\square$

Let us consider the case  $M < N$ , where the number of processors increases. We can decompose the communication matrix  $C$  in two blocks  $(A, B)$ : a left square block  $A$  of dimension  $M \times M$  and a right block  $B$  of dimension  $M \times N - M$ .

**Theorem 2.** *The communication matrix  $C = (A, B)$  is optimal if the submatrix  $A$  minimizes the migration volume and if the submatrix  $B$  minimizes the number of communication.*

*Proof.* To minimize the migration volume for  $C$ , one must take care to maximize the amount of data remaining in place, *i.e.* the sum of the terms on the diagonal of  $A$ . As a consequence,  $C$  optimizes the migration volume if  $A$  is diagonal, such as  $A = M.I_M$  with  $I_M$  the identity matrix of order  $M$ . Thus, the minimal migration volume is  $M.(N - M)$ . In this case, the number of communications of  $C$  is  $Z(C) = Z(A) + Z(B)$  with  $Z(A) = M$ . As  $B$  is assumed to be optimal,  $Z(B) = M + (N - M) - GCD(M, N - M)$  according to theorem 1. As  $GCD(M, N - M) = GCD(M, N)$ , then  $Z(C) = M + N - GCD(M, N)$  and  $C$  is optimal.  $\square$

In the case where the number of processors decreases ( $M > N$ ), we obtain a similar result by transposing the previous matrix. Theses two proofs remain correct for any perfect communication matrix, *i.e.* when  $W$  is not simply equal to  $M.N$ , but is multiple of  $M$  and  $N$ .

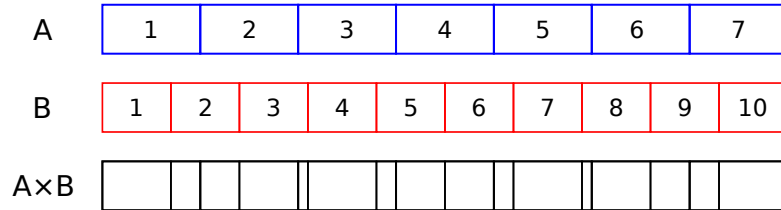


Figure 1: Partitioning of a “chain graph” (represented as a one-dimensional array of length 70) in 7 and 10 and the resulting intersection pattern  $A \times B$  used to construct the stairway communication matrix.

Let us now consider the examples given on figure 2 in the case  $7 \times 10$ . The *stairway matrix* (Fig. 2a) illustrates how to construct a perfect communication matrix with a minimum number of communications. This communication scheme is the one obtained by contiguously partitioning a “chain graph” (*i.e.* a simple one-dimensional array) in  $M$  parts and then in  $N$  parts. It is easy to demonstrate that the intersection pattern of this two partitions gives  $M + N - GCD(M, N)$  communications, which is the optimal (Fig. 1). For the stairway matrix of dimension  $7 \times 10$ , we find  $Z(C) = 16$  that is optimal, but the migration volume is not minimal at all ( $Mig(C) = 58$ ). The figure 2b gives an example of an optimal matrix, based on a stairway submatrix according to theorem 2. In this case, both the number of communication and migration volume are minimal ( $Z(C) = 16$  and  $Mig(C) = 21$ ). The figure 2c gives another example of optimal communication matrix, but not based on the stairway matrix.

In the general case where the communication matrix is not perfect (*i.e.*  $W$  is not multiple of  $M$  and  $N$ ), we can obtain similar results, defined up to an unbalance factor.

### 3.3 Correspondence Problem between the Repartitioning Hypergraph and Quotient Graph

In order to achieve a *good* repartitioning, we have to choose where to place the new parts relatively to the old ones. As the “optimal” communication scheme we want to perform during the migration phase is modeled by a repartitioning hypergraph (Sec. 3.1), we have to find a correspondence between vertices of the repartitioning hypergraph with those of the *quotient graph* associated to the initial partition (Def. 2). Indeed, vertices belonging to the same hyperedge should be matched with *close* vertices in the quotient graph as these parts will send data to the same new part.

**Definition 2.** Let  $P = (V_1, V_2, \dots, V_M)$  be the initial partition of a graph  $G$  into  $M$  parts. We note  $Q = G/P$  the quotient graph with respect to the partition  $P$ . A vertex  $i$  of  $Q$  represents the part  $V_i$  (with weight  $w(V_i)$ ) and there is an edge  $(i, j)$  in  $Q$  if the parts  $V_i$  and  $V_j$  are connected. The weight of edge  $(i, j)$  in  $Q$  is the edge-cut between parts  $V_i$  and  $V_j$ .

The *closeness* of old parts is modeled by a score. This score is computed from the edges of the quotient graph. To express this score, the repartitioning hypergraph and the quotient graph are represented by matrices. The hypergraph matrix  $H$  is a  $M \times N$  matrix and its element  $H_{v,e}$  is non-zero if the hyperedge  $e$  contains the vertex  $v$ . The quotient graph is represented by its adjacency matrix  $Q$  whose element  $Q_{i,j}$  is the weight of the edge  $(i, j)$ . A matching is represented by a  $M \times M$  permutation matrix  $X$  whose element  $X_{i,j}$  is 1 if the vertex  $i$  of  $H$  is matched with the vertex  $j$  of  $Q$ , and 0 otherwise.

In the equation 1,  $X_{i,i'}$ ,  $X_{j,j'}$ ,  $H_{i,k}$  and  $H_{j,k}$  are binary values, their product is not zero when the vertices  $i'$  and  $j'$  of  $Q$  are respectively matched with the vertices  $i$  and  $j$  of  $H$  which are in the same hyperedge  $k$ . The score is the sum of the edge weights  $Q_{i',j'}$  whose endpoints are matched with vertices belonging to the same hyperedge. Consequently, matching hyperedges of  $H$  with strongly connected subgraph of  $Q$  will give higher scores.

$$\text{score}(X) = \sum_{i,j,i',j',k} X_{i,i'} X_{j,j'} H_{i,k} H_{j,k} Q_{i',j'} \quad (1)$$

The score equation can be rewritten as follows.

$$\text{score}(X) = \sum_{i,i'} X_{i,i'} \sum_{j,j'} X_{j,j'} \left( \sum_k H_{i,k} H_{j,k} \right) Q_{i',j'} \quad (2)$$

Let  $x$  be the column vector of size  $M^2$  such that  $x_k = X_{i,i'}$  with  $k = iM + i'$  and  $\otimes$  be the Kronecker product<sup>1</sup>, the score can be rewritten as:

$$\text{score}(x) = x^T A x \quad \text{with } A = H H^T \otimes Q \text{ of size } M^2 \times M^2 \quad (3)$$

According to the previous formulation, it appears that our problem is a binary quadratic optimization problem, with linear constraints:

$$\begin{cases} \forall i, \sum_{i'} x_{iM+i'} = 1 & \text{(row constraint for } X) \\ \forall i', \sum_i x_{iM+i'} = 1 & \text{(column constraint for } X) \end{cases} \quad (4)$$

This optimisation problem is NP-hard [5]. This is a well-studied problem especially in the context of computer vision with a wide variety of applications: segmentation, clustering, graph

<sup>1</sup>Let  $A$  be a matrix of size  $P \times Q$  and  $B$  be a matrix of size  $R \times S$ . The Kronecker product  $A \otimes B$  is a block matrix with  $P \times Q$  blocks of size  $R \times S$  whose block  $(i, j)$  is  $A_{i,j} \cdot B$ .

matching. We can find in literature many heuristics to locate a good approximation to the global optimum: probabilistic metaheuristic like simulated annealing, spectral relaxation methods [9, 3], combinatorial methods like *branch & bound*, etc. In this paper, we use a basic simulated annealing with good results.

### 3.4 MxN Repartitioning Algorithm based on Fixed Vertices

Using a partitioning technique with fixed vertices in a similar way to Zoltan [2] or RM-Metis [1], the graph is repartitioned while respecting the chosen communication scheme. Our algorithm is composed of the following steps.

1. Given an initial partition  $P = (V_i)_{1 \leq i \leq M}$  of the graph  $G$  in  $M$  parts (Fig. 3a), the quotient graph  $Q$  is built (Fig. 3b).
2. An optimal communication matrices is chosen, giving us an optimal repartitioning hypergraph  $H$  (Fig. 3c). There are several possible choices as discussed in section 3.2.
3. The repartitioning hypergraph  $H$  is matched to the quotient graph  $Q$  associated with the initial partition  $P$ , using a simulated annealing algorithm to optimize the score function described in section 3.3 (Fig. 3d). It give us a permutation matrix  $X$ .
4. Fixed vertices are added to graph  $G$ . There is one fixed vertex for each new part (or hyperedge in  $H$ ). They have no weight since they represent processors, not tasks.
5. Then, we add migration edges, connected to these fixed vertices. Let  $K_j$  be the set of old processor ranks that will communicate with new processor of rank  $j$ , i.e.  $K_j = \{i \mid \exists k, X_{k,i} = 1 \wedge H_{k,j} = 1\}$ . Each fixed vertex  $j$  is connected with all the vertices of  $G$  belonging to old parts  $V_i$  with  $i \in K_j$  (Fig. 3e). These new edges are weighted with a given migration cost.
6. This enriched graph is finally partitioned in  $N$  parts, giving us the final partition  $P'$  of  $G$  (Fig. 3f).

While minimizing the edge cut, the partitioner will try to cut as few migration edges as possible, if the migration cost is high enough. Indeed, each regular vertex  $v$  of  $G$  is connected to one or more fixed vertices, modeling different new processors where  $v$  may be assigned. As exactly one of these migration edges should not be cut, the communication scheme imposed by the repartitioning hypergraph should be respected.

## 4 Experimental Results

Our  $M \times N$  graph repartitioning method is compared with a Scratch-Remap method, ParMetis and Scotch (see Sec. 2). These approaches are designed for repartitioning with a constant number of processors, but can still be used with a different new number of parts. We do not present comparison with Zoltan because it is tricky to compare graph and hypergraph cut-size, even if our approach has been generalized for hypergraph. The graph used is a simple 3D grid (of dimensions  $32 \times 32 \times 32$ ) with 32768 vertices and 95232 edges. It is initially partitioned in  $M = 8$  parts and will be repartitioned in different numbers of new parts ( $N \in [2, 24]$ ). The scratch-remap method is achieved with Scotch. It obviously provides a good cut but a high migration volume. ParMetis is used with a ratio of inter-processor communication time over redistribution time of 1000 as recommended in the documentation. Scotch remapping is used with a migration

cost of 1. The  $M \times N$  method uses a migration cost of 10 as high migration cost is needed to ensure that the communication scheme is respected.

The experiment is repeated 10 times and the charts in figure 4 show the average results. As expected from the chosen communication matrix, we can see on figure 4a that the migration for the  $M \times N$  approach is optimal. For  $N \gg M$ , the use of complex repartitioning methods becomes less relevant. The figure 4b shows that this low migration comes at the cost of higher cut, but not higher than other repartitioning tools. The cut for  $M \times N$  method is not much higher than the Scratch-Remap method which gives the better cut that the partitioner can provide with no other constraints. The number of communications (including “in-place” communications) needed for the migration is given in the figure 4c. This number is almost always optimal with our  $M \times N$  method. It can be further strengthened with an higher migration cost at the expense of edge-cut. The communication time of the migration step has been experimentally measured with OpenMPI over an InfiniBand network. The migration is up to 10% faster compared with other approaches. This confirms that our theoretical optimal communication matrices improve the migration time.

## 5 Conclusion and Future Works

We have presented in this paper a graph repartitioning algorithm, which accepts a variable number of processors, assuming the computational load is already balanced. Our algorithm minimizes both data communication and data migration overheads, while maintaining the load balance in parallel.

The preliminary experiments we have presented validate our approach comparing it against state-of-the-art partitioners. Our repartitioning provides both a minimal migration volume and a minimal number of communications, while keeping the edge cut low. Thanks to the additional constraints we provide, the results are much more stable than with other approaches.

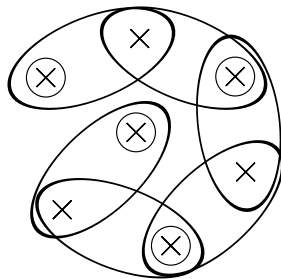
We are considering several perspectives to our work. First, we focus on graph repartitioning in the more general case where both the load and the number of processors vary. We expect this work to be really suitable for next generation of adaptive codes. Finally, to be useful in real-life applications, our algorithm needs to work in parallel and should be integrated in the *Scotch* partitioning software.

## References

- [1] Cevdet Aykanat, B. Barla Cambazoglu, Ferit Findik, and Tahsin Kurc. Adaptive decomposition and remapping algorithms for object-space-parallel direct volume rendering of unstructured grids. *J. Parallel Distrib. Comput.*, 67:77–99, January 2007.
- [2] Umit V. Catalyurek, Erik G. Boman, Karen D. Devine, Doruk Bozdağ, Robert T. Heaphy, and Lee Ann Riesen. A repartitioning hypergraph model for dynamic load balancing. *J. Parallel Distrib. Comput.*, 69(8):711–724, 2009.
- [3] Olivier Duchenne, Francis R. Bach, In-So Kweon, and Jean Ponce. A tensor-based algorithm for high-order graph matching. In *CVPR*, pages 1980–1987. IEEE, 2009.
- [4] S. Fourestier and F. Pellegrini. Adaptation au repartitionnement de graphes d’une méthode d’optimisation globale par diffusion. In *Proc. RenPar’20, Saint-Malo, France*, May 2011.
- [5] Michael R. Garey and David S. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman, 1979.
- [6] Bruce Hendrickson and Karen Devine. Dynamic load balancing in computational mechanics. In *Computer Methods in Applied Mechanics and Engineering*, volume 184, pages 485–500, 2000.
- [7] Bruce Hendrickson, Robert W. Leland, and Rafael Van Driessche. Skewed graph partitioning. In *Eighth SIAM Conf. Parallel Processing for Scientific Computing*, 1997.
- [8] Saeed Iqbal and Graham F. Carey. Performance analysis of dynamic load balancing algorithms with variable number of processors. *Journal of Parallel and Distributed Computing*, 65(8):934 – 948, 2005.
- [9] Marius Leordeanu and Martial Hebert. A spectral technique for correspondence problems using pairwise constraints. In *Proceedings of the Tenth IEEE International Conference on Computer Vision - Volume 2, ICCV ’05*, pages 1482–1489, Washington, DC, USA, 2005. IEEE Computer Society.
- [10] Leonid Oliker and Rupak Biswas. Plum: parallel load balancing for adaptive unstructured meshes. *J. Parallel Distrib. Comput.*, 52:150–177, August 1998.
- [11] J.R. Pilkington and S.B. Baden. Dynamic partitioning of non-uniform structured workloads with spacefilling curves. *Parallel and Distributed Systems, IEEE Transactions on*, 7(3):288–300, March 1996.
- [12] Kirk Schloegel, George Karypis, and Vipin Kumar. Multilevel diffusion schemes for repartitioning of adaptive meshes. *Journal of Parallel and Distributed Computing*, 47(2):109 – 124, 1997.
- [13] James D. Teresco, Karen D. Devine, and Joseph E. Flaherty. Partitioning and dynamic load balancing for the numerical solution of partial differential equations. In Timothy J. Barth, Michael Griebel, David E. Keyes, Risto M. Nieminen, Dirk Roose, Tamar Schlick, Are Magnus Bruaset, and Aslak Tveito, editors, *Numerical Solution of Partial Differential Equations on Parallel Computers*, volume 51 of *Lecture Notes in Computational Science and Engineering*, pages 55–88. Springer Berlin Heidelberg, 2006.

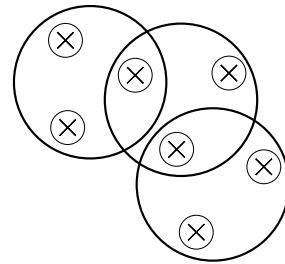
- [14] Rafael Van Driessche and Dirk Roose. Dynamic load balancing with a spectral bisection algorithm for the constrained graph partitioning problem. In Bob Hertzberger and Giuseppe Serazzi, editors, *High-Performance Computing and Networking*, volume 919 of *Lecture Notes in Computer Science*, pages 392–397. Springer Berlin / Heidelberg, 1995.

$$\begin{bmatrix} 7 & 3 & & & & & & & & & \\ & 4 & 6 & & & & & & & & \\ & & 1 & 7 & 2 & & & & & & \\ & & & 5 & 5 & & & & & & \\ & & & & 2 & 7 & 1 & & & & \\ & & & & & 6 & 4 & & & & \\ & & & & & & 3 & 7 & & & \end{bmatrix}$$



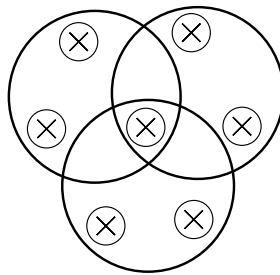
(a) Stairway matrix.

$$\begin{bmatrix} 7 & & & & & & & & & & 3 \\ & 7 & & & & & & & & & 3 \\ & & 7 & & & & & & & & 1 & 2 \\ & & & 7 & & & & & & & 3 \\ & & & & 7 & & & & & & 2 & 1 \\ & & & & & 7 & & & & & 3 \\ & & & & & & 7 & & & & 3 \\ & & & & & & & 7 & & & 3 \end{bmatrix}$$



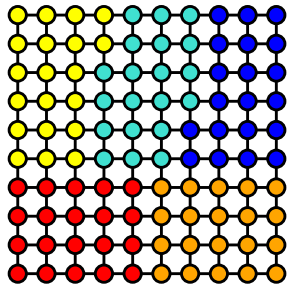
(b) Optimal matrix based on stairway submatrix.

$$\begin{bmatrix} 7 & & & & & & & & & & 3 \\ & 7 & & & & & & & & & 3 \\ & & 7 & & & & & & & & 3 \\ & & & 7 & & & & & & & 3 \\ & & & & 7 & & & & & & 3 \\ & & & & & 7 & & & & & 3 \\ & & & & & & 7 & & & & 3 \\ & & & & & & & 7 & & & 3 \\ & & & & & & & & 7 & & 3 \\ & & & & & & & & & 7 & 1 & 1 & 1 \end{bmatrix}$$

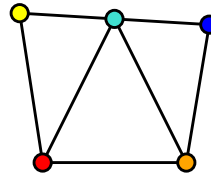


(c) Another optimal matrix.

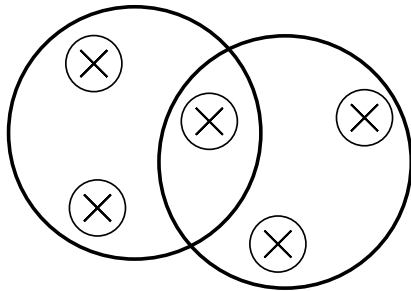
Figure 2: Three communication matrices in the case  $7 \times 10$  and their representation as repartitioning hypergraph. Zero elements in matrices are not shown. The elements in red are those who remain in place during communications, others will migrate.



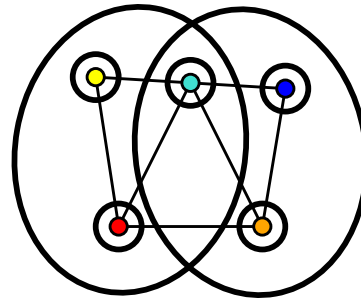
(a) Initial partition in 5 parts.



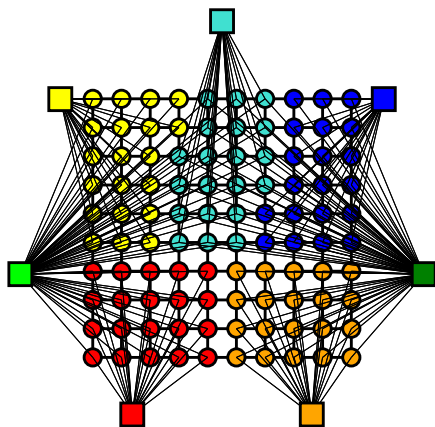
(b) Quotient graph of the initial partition.



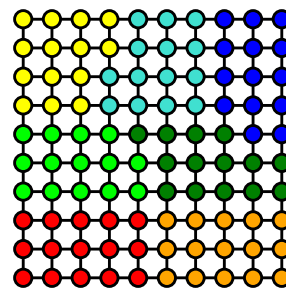
(c) An optimal repartitioning hypergraph for the case  $5 \times 7$ .



(d) Matching between the quotient graph and the repartitioning hypergraph.



(e) Graph with fixed vertices added according to the matching.



(f) Final partition in 7 parts.

Figure 3: Repartitioning from 5 to 7 parts.



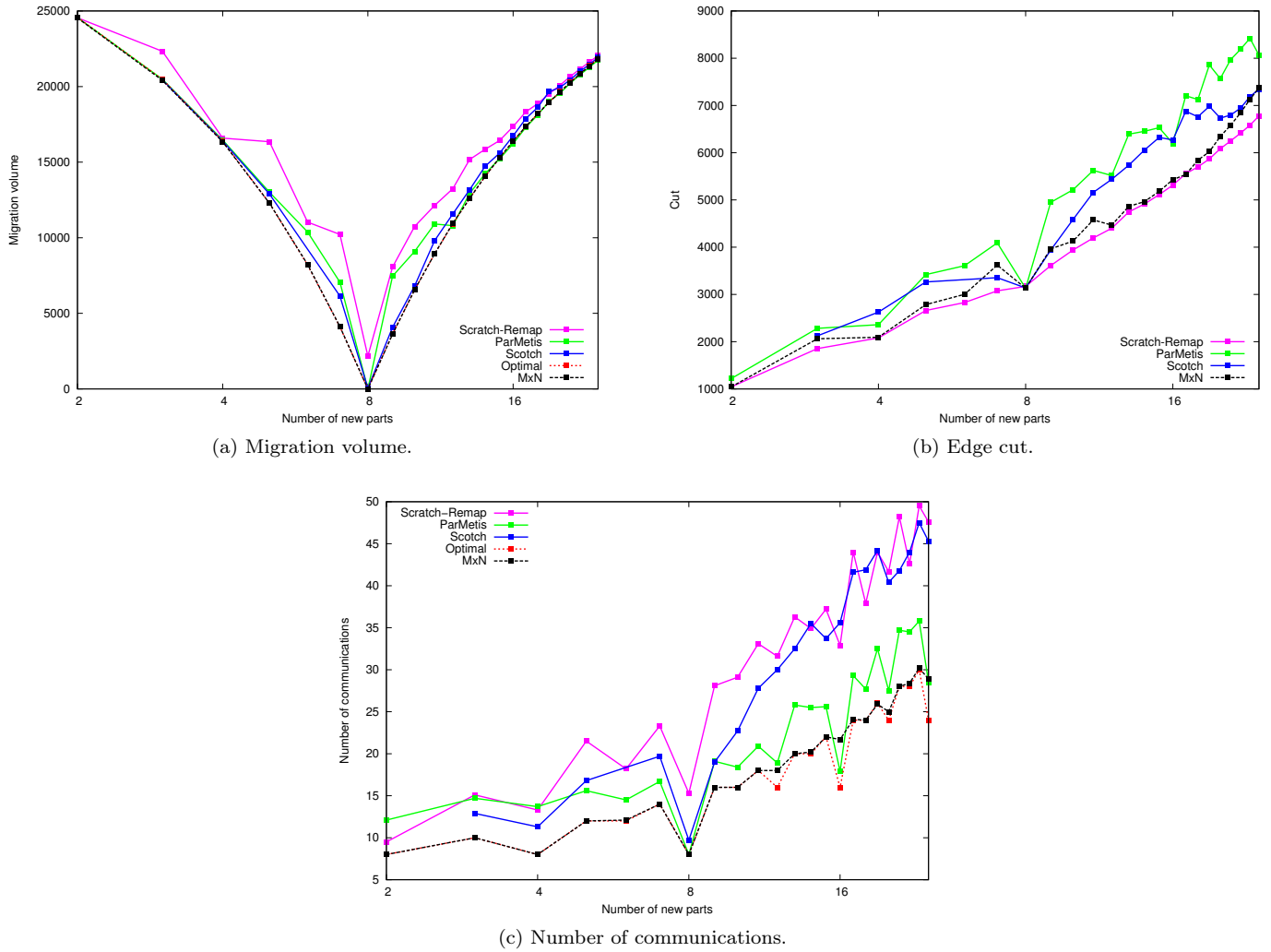


Figure 4: Experimental results for the graph repartitioning of a  $32 \times 32 \times 32$  grid from  $M = 8$  processors to  $N \in [2, 24]$ .



**RESEARCH CENTRE  
BORDEAUX – SUD-OUEST**

351, Cours de la Libération  
Bâtiment A 29  
33405 Talence Cedex

Publisher  
Inria  
Domaine de Voluceau - Rocquencourt  
BP 105 - 78153 Le Chesnay Cedex  
[inria.fr](http://inria.fr)

ISSN 0249-6399