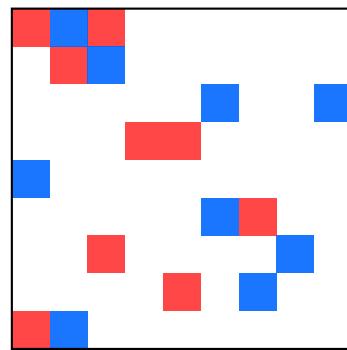


Iterative sparse matrix partitioning



Author:
Davide Taviani

Supervisor:
Prof. dr. Rob H. Bisseling
Second reader:
Dr. Paul A. Zegeling

Master Thesis
October 7, 2013



Universiteit Utrecht

Preface

This master thesis is, technically, the fruit of my labor over the past nine months, but I would say it is more the result of these two years here at Utrecht University. When I first started, “scientific computing” was still a somewhat cloudy concept; when I discovered the field of combinatorial optimization, a new world opened before my eyes: finally my mathematical thinking could be used to solve a number of problems which lie at the very base of computer science, my other great passion. To me, the most amazing aspect of this field is how the basic principles and concepts are surprisingly easy to understand, and yet finding the best solutions is all but trivial.

I think the simplicity of the problems is well reflected also in my thesis: the goal is efficient sparse matrix partitioning, in order to have a low communication volume during sparse matrix-vector multiplications. Although it might sound complicated, the concept is extremely easy to grasp, even without any familiarity with mathematics: given a very large grid of numbers (matrix), which is mostly empty (sparse), you want to color the nonempty cells either red or blue, such that, at the end, there is roughly the same number of red cells than blue cells; now, the goal is to find a way of doing so such that, in our matrix, rows and columns have, *as much as possible*, only one color. The simplicity of the subject is also well reflected, in my opinion, by the large amount of pictures in this thesis: whenever possible, I tried to come up with simple visualizations that ease the understanding of the discussed concepts.

Acknowledgments

First of all, I would like thank my supervisor Rob Bisseling: his guidance was fundamental during these months of research and his suggestions extremely useful. Moreover, I want also to thank him for being “responsible” for most of my excitement toward the field. Secondly, I want to thank Bas Fagginger Auer, whose help was vital: every time I was stuck with a problem (mostly with code, but not exclusively) he eagerly, and effectively, helped me with it. I also want to thank him for teaching me a lot, during my first year, in the topic of mathematical writing and this thesis was heavily influenced by him.

In aggiunta, vorrei ringraziare anche la mamma e il papà: è al loro aiuto e supporto che devo davvero tutto. Un grazie anche agli amici, ché durante questi mesi mi hanno sempre incoraggiato e hanno dimostrato sempre molta pazienza nei miei confronti. Infine, un grande grazie a Kleopatra, il cui aiuto non poteva essere più prezioso: grazie per avermi sopportato e supportato anche nei momenti più difficili di questi mesi di ricerca, aiutandomi sempre a trovare lo spirito giusto per affrontare i problemi.

Apart from the friends and family back home (I hope the reader will forgive me for having momentarily switched to Italian), I also want to thank my friends here in Utrecht, who helped me greatly to adapt to the Dutch life and were very useful in helping me keeping stress levels down.

At the end, I would like also to thank you, reader. If you are reading this line after the others, you at least read one page of my thesis, which already fills me with joy. Thank you.

Contents

Preface	I
1 Introduction	1
1.1 Parallel sparse matrix-vector multiplication	1
1.2 Hypergraph model	5
1.3 Earlier work	8
1.4 Medium-grain method	9
2 Strategies for splitting a matrix A into A_r and A_c	11
2.1 Individual assignment of nonzeros	11
2.2 Assignment of blocks of nonzeros	12
2.2.1 Using the Separated Block Diagonal form of A	13
2.2.2 Using the Separated Block Diagonal form of order 2 of A	15
2.3 Maximizing empty rows of B	18
2.4 Partial assignment of rows and columns	19
2.4.1 Computation of the priority vector v	20
3 Maximum independent set formulation of the partial row/column assignment problem	24
3.1 Graph construction	24
3.2 The maximum independent set and its computation	25
3.2.1 Maximum independent set	25
3.2.2 The Hopcroft-Karp algorithm for bipartite matching	28
3.3 Computation of the priority vector v with the maximum independent set	29
4 Implementation and experimental results	31
4.1 Test matrices	32
4.2 Preliminary selection of the best heuristics	36
4.2.1 Partition-oblivious heuristics	36
4.2.2 Partition-aware heuristics	38
4.3 Number of iterations	40

4.4	Analysis of the performance of the best heuristics	44
4.4.1	Partition-oblivious heuristics	44
4.4.2	Partition-aware heuristics	45
4.5	Results with iterative refinement of the final partitionings	45
5	Conclusions and further developments	50
5.1	Improvement of the initial solution	50
5.2	Iterative partitioning	50
5.3	Further research	51
	Bibliography	52

Chapter 1

Introduction

1.1 Parallel sparse matrix-vector multiplication

Matrices are one of the most important mathematical objects, as they can be used to represent a wide variety of data in many scientific disciplines: they can encode the structure of a graph, define Markov chains with finitely many states, or possibly represent linear combinations of quantum states or also the behaviour of electronic components.

In most real-world computations, the systems considered are usually of large size and involve **sparse** matrices, because the variables at hand are usually connected to a limited number of others (for example, a large graph in which each node has just a handful of incident edges); therefore, the matrices involved have the vast majority of entries equal to 0. More formally, let us consider a matrix of size $m \times n$ with N nonzeros. We say that the matrix is sparse if $N \ll mn$. Without loss of generality, we assume that each row and column has at least one nonzero (otherwise those rows and columns can easily be removed from the problem).

One of the most fundamental operations performed in these real-world computations is the sparse matrix-vector multiplication, in which we compute

$$u := Av, \tag{1.1}$$

where A denotes our $m \times n$ sparse matrix, v denotes a dense vector of length n , and u the resulting vector of length m .

The computation of this quantity following the definition of matrix-vector multiplication, i.e. with the sum

$$u_i = \sum_{j=0}^{n-1} a_{ij} v_j, \quad \text{for } 0 \leq i < m,$$

requires $\mathcal{O}(mn)$ operations; this is not efficient with a sparse matrix: if we perform the multiplications only on the nonzero elements, we obtain an algorithm with running time $\mathcal{O}(N)$, and by definition of sparsity, $N \ll mn$.

As mentioned, the systems considered are large, with sparse matrices with thousands (even millions) of rows and columns and millions of nonzeros; for such big instances, even a running time of $\mathcal{O}(N)$ might be non-negligible, especially since sparse matrix-vector multiplications are usually just a part of a bigger iterative algorithm, and need to be performed several times.

It is an important goal then to be able to perform such computations in the least amount of time possible:

however, as there is a natural tradeoff between power consumption and the speed of the processing units [1], it is not feasible to rely only on fast CPUs, but rather focus on parallelism and employ a large number of them with lower processing speed (and, as a result, with fairly low energy requirements).

To describe an efficient way of performing parallel sparse matrix-vector multiplications, we follow the approach described in [2]: before the actual computation takes place, the sparse matrix is distributed among the p processors, creating a **partitioning** of the set of the nonzeros: A is split into disjoint subsets A_0, \dots, A_{p-1} . Moreover, also the input vector v and the output vector u are distributed among the p processors (note that their distribution might not necessarily be the same, and often it is not).

Figure 1.1 shows a possible partitioning of a 9×9 matrix with 18 nonzeros. As the actual values of the nonzeros are not important, we only show the sparsity pattern (a colored cell means that there is a nonzero in that position). The two colors denote, respectively, the two resulting subsets of nonzeros.

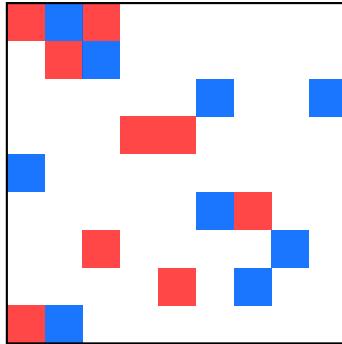


Figure 1.1: Example of a distribution among two processors of a 9×9 matrix with 18 nonzeros. Only the sparsity pattern is shown.

After this distribution, every processor has to compute its local contribution toward the matrix-vector multiplication: to do so, it requires the appropriate vector components which might have been assigned to another processor during the data distribution; if this is the case, communication is required. Once all the required vector components are obtained, the processor starts computing all its local contributions, which are afterwards sent to their appropriate owner, according to the distribution of u . The three phases that describe this process for processor $s = 0, \dots, p - 1$, are summarized in Algorithm 1.1, from [2, 3].

In reality there is also a fourth phase, in which each processor sums up all the contributions received in phase (2) for all of its owned components of u ; this is a small phase with negligible computational cost compared to phase (2) and for this reason it has been omitted from the algorithm.

Note that we assume that all of the nonzero values are represented with the same amount of bits. Doing so, we can focus exclusively on the coordinates of the nonzeros, omitting completely their values, as it does not influence the cost of a parallel sparse matrix-vector multiplication.

Figure 1.2 shows an example of the communication involved in phases, or supersteps (0) and (2), for the example partitioning shown in Figure 1.1: the vertical arrows represent the fan-out, while the horizontal arrows represent the fan-in; the color of an arrow indicates which processor is sending data.

As our main interest is to **minimize** the time spent by the parallel machine computing this sparse matrix-vector multiplication, we need to compute explicitly the cost of Algorithm 1.1: we can immediately note that such an algorithm, which follows the Bulk Synchronous Parallel model [4], consists of two communication supersteps separated by a computation superstep.

The time spent (in the sense of number of operations performed) by a parallel machine in a computation superstep is exactly the time taken by the processor that finishes last: more formally, the time cost of step (1):

Input: A_s , the local part of the vector v

Output: The local part of the vector u

```

 $I_s := \{i | \exists j : a_{ij} \in A_s\}$ 
 $J_s := \{j | \exists i : a_{ij} \in A_s\}$ 

(0) ▷ Fan-out
for all  $j \in J_s$  do
    Get  $v_j$  from the processor that owns it.
end for

(1) ▷ Local sparse matrix-vector multiplication
for all  $i \in I_s$  do
     $u_{is} := 0.$ 
    for all  $j$  such that  $a_{ij} \in A_s$  do
         $u_{is} = u_{is} + a_{ij}v_j.$ 
    end for
end for

(2) ▷ Fan-in
for all  $i \in I_s$  do
    Send  $u_{is}$  to the owner of  $u_i$ .
end for

```

Algorithm 1.1: Parallel sparse matrix-vector multiplication.

$$T_{(1)} = \max_{0 \leq s < p} |A_s|. \quad (1.2)$$

where $|A_s|$ denotes the size of A_s , the part of A assigned to processor s .

It is easy to understand that, in order to have efficient parallelization, the computation load has to be distributed evenly. Usually, however, it is not possible to achieve a perfect load balance (e.g. when dividing an odd number of computations among an even number of processors) and we have to reason in terms of an allowed imbalance ε . Consequently, we impose the following hard constraint about the maximum size of the subsets of nonzeros assigned to each processor, according to [2, eq. 4.27]:

$$\max_{0 \leq s < p} |A_s| \leq (1 + \varepsilon) \frac{N}{p}. \quad (1.3)$$

Typical values for the allowed ε in this constraint are 0.03, i.e. a 3% imbalance.

It is reasonable, after all, that the problem of finding an efficient way of performing this computation step boils simply down to a hard constraint for the data distribution. This is because we still have to perform all the multiplications of the form $a_{ij}v_j$, no matter our choice. The communication costs, represented by the first and last supersteps in Algorithm 1.1, are instead the most interesting aspect about maximizing the efficiency of a parallel sparse matrix-vector multiplication algorithm, as there is extreme variability. As a simple example, suppose $p = 2$ and consider the matrix represented in Figure 1.3.

Two possible partitionings of this matrix into two sets are given in Figure 1.4. In Figure 1.4(a) no communication is necessary, whereas in Figure 1.4(b), all of the rows and columns are split, and therefore the maximum possible communication is required during the sparse matrix-vector multiplication algorithm.

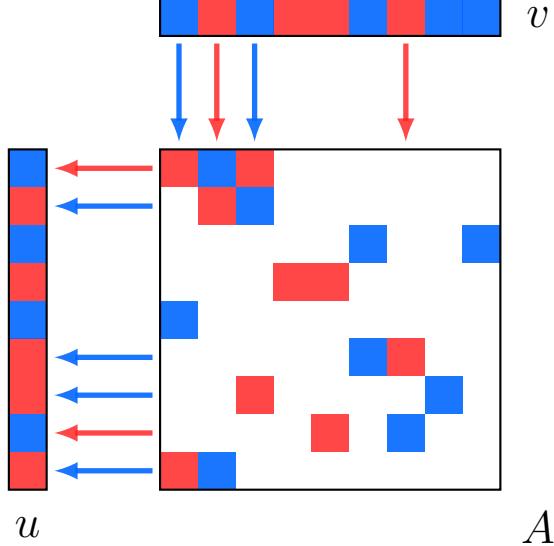


Figure 1.2: Communication for the parallel sparse matrix-vector multiplication with a matrix partitioned as in Figure 1.1. Vertical arrows represent step (0) while horizontal ones represent step (2). The color of an arrow denotes which processor is sending their data for that row/column.

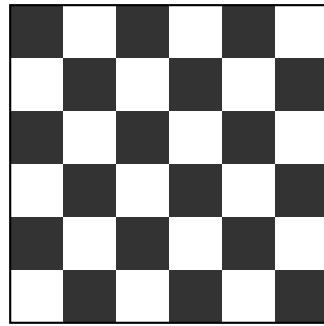


Figure 1.3: Example matrix with checkered sparsity pattern. Black boxes represent the nonzeros.

Previously, we claimed that the matrix and both the vectors have to be partitioned: in reality it is sufficient to consider only the problem of distributing the nonzeros, and the partitioning of the vector can be executed according to this: because of the structure of the communication supersteps in Algorithm 1.1, communication is required if and only if the rows/columns of the matrices are *cut*, i.e. assigned to more than one processor.

If a full column of our matrix A is assigned to the same processor, we can freely assign the corresponding component of v to the same processor, eliminating completely one source of communication (namely, the fan-out for that column). The same reasoning can be done for the rows. This simplification is possible because imposing a hard constraint similar to (1.3) also to the vector distribution is not helpful, as it only affects the time of linear vector operations outside the matrix-vector multiplication, which are in generally much cheaper [3, Sec. 3].

We can describe more formally the communication costs, following the notation of [3, Def. 2.1]: let A_0, \dots, A_{p-1} be a p -way (with $p \geq 1$) partitioning of the sparse matrix A of size $m \times n$. Let λ_i denote the number of processors which have a nonzero of row i and let μ_j be the number of processors that have a nonzero of column j ; note that, because we assumed that all the rows and columns are nonempty, $\lambda_i, \mu_j \geq 1$. Then the total time costs for the communication steps in our Algorithm 1.1 are:

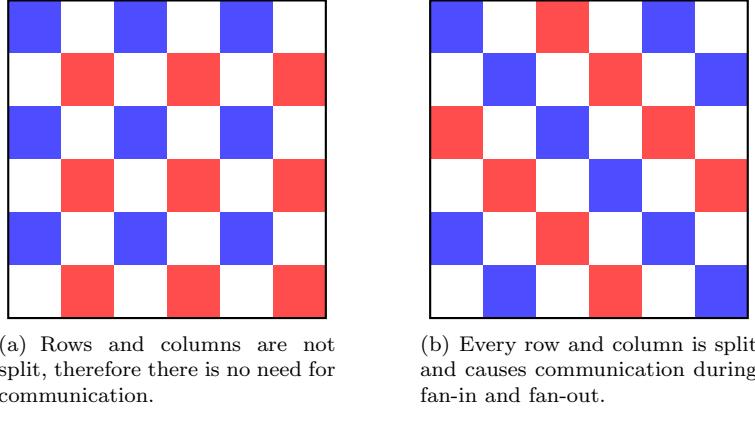


Figure 1.4: Different partitionings of the matrix from Figure 1.3. Red and blue squares represent nonzeros assigned to the two different processors.

$$\begin{aligned} T_{(0)} &= \sum_{j=0}^{n-1} (\mu_j - 1), \\ T_{(2)} &= \sum_{i=0}^{m-1} (\lambda_i - 1). \end{aligned} \tag{1.4}$$

These costs are quite straightforward: it is reasonable to assume that the owner of the appropriate vector component is one of the processors that have a nonzero in that row/column, and therefore communication is not necessary for that processor. Adding these costs together, we define the **communication volume** V of the considered partitioning as

$$V := V(A_0, \dots, A_{p-1}) = T_{(0)} + T_{(2)} = \sum_{i=0}^{m-1} (\lambda_i - 1) + \sum_{j=0}^{n-1} (\mu_j - 1). \tag{1.5}$$

As we can see, the communication volume V depends entirely on the matrix A and the considered partitioning. Therefore, the problem of minimizing the cost of a matrix-vector multiplication is shifted toward finding an efficient way of distributing the sparse matrix among the available processors, such that our balance constraint (1.3) is satisfied. The following sections and chapters and, ultimately, this whole Master Thesis, are therefore dedicated to it.

1.2 Hypergraph model

The problem of distributing the nonzeros of a matrix in order to minimize the communication volume, or, in short, the matrix partitioning problem, can also be viewed from the graph theory point of view. We recall that a (unweighted, undirected) graph $G = (V, E)$ is a set of vertices (or nodes) V and edges E which connect them.

The graph partitioning problem has been used in the past to model the load balancing in parallel computing: data are represented as vertices, while their connections (the dependencies) are represented with edges. For a more rigorous definition of the graph partitioning problem, we follow the notation given in [5], performing the simplification in which all the edges have unitary weight. Given the graph $G = (V, E)$ we say that (V_0, \dots, V_{p-1}) is a p -way partition of G if all these subsets are nonempty, mutually disjoint and their union is the whole set of nodes V .

Moreover, we can consider a balance criterion similar to (1.3):

$$\max_{0 \leq s < p} |V_s| \leq (1 + \varepsilon) \frac{|V|}{p}, \quad (1.6)$$

where ε , similarly as before, represents the allowed imbalance.

Now, given a partition (V_0, \dots, V_{p-1}) of the graph G , we say that the edge $e = (i, j)$ is *cut* if $i \in V_k, j \in V_l$, with $k \neq l$; otherwise, it is said to be *uncut*. Previously, we claimed that communication during the parallel matrix-vector multiplication can be avoided if a row/column is uncut, and here the goal is the same: we want to minimize the *cutsize*, i.e. the number of edges cut.

However, despite all the similarities between the matrix partitioning problem and the graph partitioning one, it has been shown [5] that this cut-edge metric is not an accurate representation of the communication volume. Additional criticism [6] comes from the fact that the graph partitioning approach can only handle square symmetric matrices, even though it has been shown to be a good enough approximation in that case [7]. Moreover, it was also shown [8] that these disadvantages hold for all application of graph partitioning in parallel computing, and not only our problem of matrix partitioning for sparse matrix-vector multiplication. An exact way of modeling the matrix partitioning problem is through the concept of hypergraph partitioning [5].

A hypergraph is simply a generalization of a graph: we do not consider edges that connect two nodes, but rather *hyperedges* (or *nets*), which are subsets of nodes. Apart from considering only non-empty hyperedges, note that there is no other restriction on their cardinality.

Hypergraphs, and in particular the hypergraph partitioning problem are already well known in literature: they have a natural application in the design of integrated circuits (VLSI), in finding efficient storage of large databases on disks, and data mining [9], as well as urban transportation design and study of propositional logic [10].

Because of this extensive application basis, translating our matrix partitioning problem to a hypergraph partitioning problem seems quite convenient, as all the methods already developed can be analyzed and employed also in our case.

Figure 1.5 shows an example of a hypergraph. Each colored set represents a different hyperedge; we can see that we can have hyperedges which contain only one node.

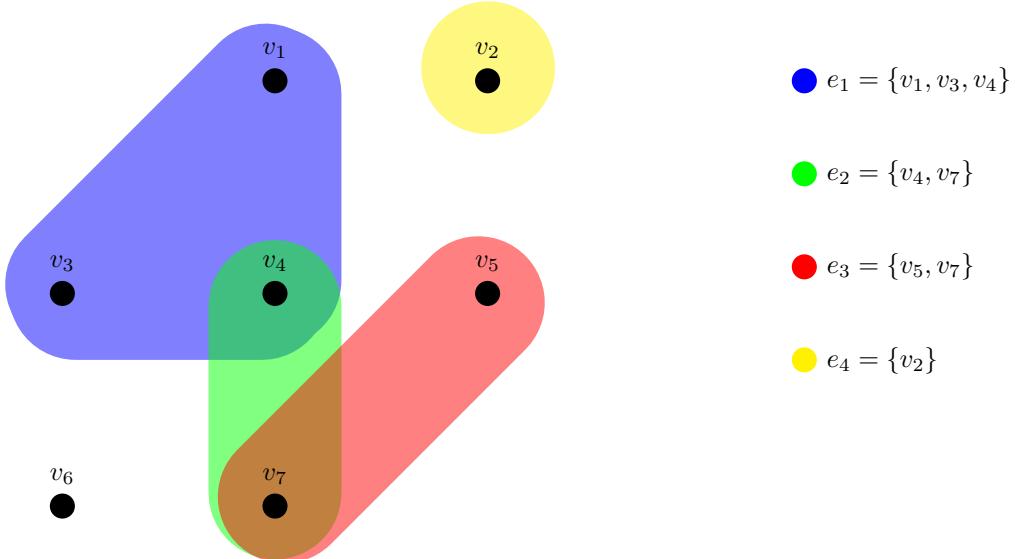


Figure 1.5: Example of a hypergraph with 7 nodes and 4 hyperedges.

The definition of hypergraph partitioning problems is identical to the case of a graph, with the difference that now we do not have cut edges, but cut hyperedges: given the hyperedge $e = \{v_1, \dots, v_k\}$, we say

that e is cut if there are i, j such that $v_i \in V_r, v_j \in V_s$, with $r \neq s$, i.e. at least two nodes belong to different sets of the partition. As usually, we want to minimize the cut hyperedges.

If, similarly to (1.2), we define λ_e as the number of different sets the vertices in the hyperedge e are assigned to, the total cost of the partition (V_0, \dots, V_p) is:

$$C = C(V_0, \dots, V_{p-1}) = \sum_{e \in E} (\lambda_e - 1). \quad (1.7)$$

We can see how closely these equations resemble the ones given in the previous section: it is clear that the hypergraph partitioning problem closely resembles our original matrix partitioning problem.

Note that the hypergraph partitioning problem, along with the graph partitioning problem, are known to be NP-hard [11, Ch. 6].

Now, we will describe three possible models for the decomposition of a sparse matrix into a hypergraph, and discuss their advantages and disadvantages.

In the **column-net** model, our matrix A is represented as a hypergraph for a row-wise decomposition: rows of the matrix are nodes ($V = \{v_1, \dots, v_m\}$, while columns are hyperedges ($E = \{e_1, \dots, e_n\}$). We have that the node v_i belongs to the hyperedge e_j (in short $v_i \in e_j$) if and only if $a_{ij} \neq 0$. With this model, the size of the hyperedge e_j is exactly the number of nonzeros in column j , whereas the node v_i belongs exactly to as many hyperedges as there are nonzeros in row i .

As already said, performing a partitioning on the hypergraph consists of assigning each vertex to one of the sets V_0, \dots, V_{p-1} . In this model, this corresponds to assigning a row completely to a processor.

However, as vertices are not exactly nonzeros of our matrix, (1.3) and (1.6) are not exactly equivalent; we need to adjust our balance constraint by introducing a weight for each vertex, as in [2, Def. 4.34]. For $v_i \in V$, we define its weight c_i as

$$c_i := |\{j : a_{ij} \neq 0\}|,$$

which simply is the number of nonzeros in row i of the matrix A . Note that, following the same notation as in the previous section, we can see the total number of nonzeros N as $N = \sum_{v_i \in V} c_i$.

Our modified balance constraint is as follows:

$$\max_{0 \leq s < p} W(V_s) := \max_{0 \leq s < p} \sum_{v_i \in V_s} c_i \leq (1 + \varepsilon) \frac{N}{p}. \quad (1.8)$$

The **row-net** model is similar to the one just described (as can be guessed from the name): it is exactly the transposed of the column-net model, in the sense that now rows are hyperedges and columns are vertices of the hypergraph. The reasoning just described applies also to this model, with the little modification that now the weight of a vertex is the number of nonzeros in that column.

We see how the column-net model and row-net model have the advantage of fully assigning a row (or a column) to a processor; this has the advantage of eliminating completely one source of communication in our parallel sparse matrix-vector multiplication algorithm (respectively, the fan-in and fan-out). However, this advantage can easily become a weakness, because now the partitioning is forcedly 1-dimensional, and this is usually too strong of a restriction.

Now, as a last example of possible decomposition of a matrix into a hypergraph, and as a partial address to the drawbacks of the previous two models, we will describe a 2-dimensional approach, the so-called *fine-grain* model [12]. In this model, the N nonzeros are the vertices ($V = \{v_1, \dots, v_N\}$) and the m rows and n columns are hyperedges ($E = E_r \cup E_c = \{e_1, \dots, e_m\} \cup \{e_{m+1}, \dots, e_{m+n}\}$). With this notation, E_r represents the row hyperedges and E_c represents the column hyperedges. The relationship between the vertices and the hyperedges is fairly obvious: $v_k = a_{ij}$ is in both e_i and e_{m+j} .

Now, as the vertices correspond exactly to nonzeros of our matrix, we can use the original equation (1.6) as balance constraint; if we combine this with (1.7), which describes the cost of a hypergraph partition, we can clearly see how this is identical to our original matrix partitioning problem, described by (1.3) and (1.5).

On a higher level, one of the benefits of this decomposition model is easy to understand: we have a lot of freedom and we can assign individually each nonzero to a different partition. Similarly as before, however, this advantage can easily become a drawback because now the size of the hypergraph is consistently larger, with N vertices compared to m and n of the previous two models. Thus, computations on the fine-grain model take substantially more time than row-net or column-net models and therefore there is a restriction on the size of the problem that can be efficiently solved.

1.3 Earlier work

Among the models used to translate matrix partitioning into hypergraph partitioning, we already mentioned row-net and column-net [5], proposed in 1999, and a more recent fine-grained approach [12], proposed in 2001.

In addition to these models, there has been some research effort towards the creation of more sofisticated methods, which often comprise several stages and combine different models. Recently Pelt and Bisseling proposed the **medium-grain** method [13]. As this method is at the base of our work, a more detailed explanation will be given in Section 1.4.

For example, Uçar and Aykanat [14] first employ an elementary 1-dimensional hypergraph model, and then they transform it in several ways to different hypergraph models suitable for both symmetric and unsymmetric matrix partitionings; it is important to note that these models also include the input and output vectors, and therefore a few extra vertices are added to the hypergraph.

A different 2-dimensional approach is given by the *coarse-grain* method [15]: first the column-net hypergraph model is used, obtaining a row partitioning of the matrix in p parts, then a multi-constraint column partitioning in q parts is performed, yielding a final 2-D cartesian partitioning in $p \times q$ parts.

Moreover, Vastenhoud and Bisseling proposed a 2-dimensional recursive method for data distribution [3]; this greedy method splits recursively a rectangular matrix into 2 parts. At each step of the recursion, there is the choice on the direction to be taken in the next step: two different strategies are proposed, alternating splitting directions or simply trying to split both vertically and horizontally and taking greedily the best of the two.

Besides these general purpose models and methods, it is also possible to take into account the structure of the matrix to be partitioned: Hu, Maguire, and Blake present in [16] an algorithm for nonsymmetric matrices that performs row and column permutations, obtaining a bordered block diagonal form and then trying to assign matrix rows such that the number of cut columns is minimized.

In general, as there is such a wide variety of different methods and models, it might be difficult to choose the best one, given a matrix to partition. Çatalyürek, Aykanat, and Uçar propose a partitioning recipe [17] that chooses a partitioning method according to some matrix characteristics.

Regarding the actual implementations of the just discussed models, methods and algorithms, there are a few existing software partitioners available. Among the sequential ones we have PaToH (a multilevel Partitioning Tool for Hypergraphs) [18], hMetis [19] (specifically targeted at partitioning hypergraphs for VLSI design), Mondriaan [3] (among the ones here described, this is the one more specifically designed to solve the matrix partitioning problem), MONET (Matrix Ordering for minimal NET-cut)[16]. Zoltan-PHG (Parallel Hypergraph Partitioner) [20] performs instead matrix partitioning in parallel; the relative scarcity of parallel software partitioners is to be explained by the fact that this field is relatively new, and therefore most of the research efforts have been directed toward a sequential approach.

The partitioners just mentioned produce close but slightly different results, with respect to both solution quality and execution time, having at the core the same method for finding good initial solutions. The initial partitioning method employed for small subproblems is the well-known Kernighan-Lin [21] method,

with the optimizations of Fiduccia-Mattheyses [22]. This local search heuristic was originally designed for bipartitioning graphs and, given a partitioning that obeys the balance constraint (1.3), it applies a series of small changes to improve the quality of the solution.

To solve large instances, all these partitioners use a multi-level method: the large problem is progressively coarsened until a smaller instance is obtained, then the problem is solved on this small instance and the solution is gradually uncoarsened, with a refinement at each step to improve the solution quality.

Finally, these existing software partitioners are all based on *recursive bisection*: instead of partitioning the hypergraph directly into the desired number of parts, they execute a sequence of bisections of the partitions. This is a good simplification in the sense that it just suffices to find good algorithms for bipartitioning, and also because splitting a hypergraph in just two parts is much easier; there is however one major flaw with this approach: using this recursive bisection we might not be able to reach the same quality of a solution as with direct splitting into the desired number of parts. This flaw is, however, often ignored and direct splitters into more than two parts are rare.

1.4 Medium-grain method

All of the possible ways of translating the matrix partitioning problem into a hypergraph partitioning problem have different advantages and drawbacks: the 1-dimensional ones, row-net and column-net, eliminate completely one source of communication but are somewhat too restrictive; fine-grain, on the contrary, does not provide any kind of limitation on the choices for the partitioning, but the resulting hypergraph is often too big to manage.

A new method has recently been proposed by Pelt and Bisseling [13], which can be described as a sort of middle ground between the 1-dimensional models and fine-grain model. The resulting partitioning is 2-dimensional by design (thus avoiding the limitations of the row-net and column-net models), but it still imposes that clusters of nonzeros from the same rows and columns are assigned to the same processor, thus reducing the size of the final hypergraph, avoiding the main disadvantage of the fine-grain model.

The key of the medium-grain method lies into the splitting of our original matrix A in two parts, A_r and A_c , such that $A_r + A_c = A$. Then, we proceed to construct the auxiliary block-matrix B , of size $(m+n) \times (m+n)$, defined as

$$B := \begin{bmatrix} I_n & A_r^T \\ A_c & I_m \end{bmatrix}, \quad (1.9)$$

where I_n and I_m denote, respectively, the identity matrices of size n and m . The final hypergraph is finally obtained by applying the row-net model to this matrix B .

Figure 1.6 illustrates this process for a 3×6 rectangular matrix A .

After we apply the row-net model and obtain a partitioning of the hypergraph, it is immediate to retrieve a partitioning of our matrix A , as depicted in Figure 1.7.

The usefulness of A_c and A_r is clear if we consider that we use the row-net model. The first is left as-is, while the second is transposed; then, when partitioning 1-dimensionally such that the columns are kept together, we see that we are effectively keeping together elements within the same columns of A_c and A_r^T . The resulting partitioning is fully 2-dimensional, because there are clusters of nonzeros: rows for A_r and columns for A_c (hence the subscripts).

The diagonal elements of B are used only to compute the communication volume. Let us consider the k th column of A ; the corresponding nonzeros can be found in the k th column of A_c and in the k th row of A_r^T . If both these parts are nonempty, i.e. the k th column of A was not fully assigned to either A_r or A_c , we need to be careful when we compute the communication volume of a given partitioning: if these parts are to different processors, communication is needed in Algorithm 1.1.

Therefore the diagonal nonzero $B_{k,k}$, assigned by the row-net model to the same processor as the k th

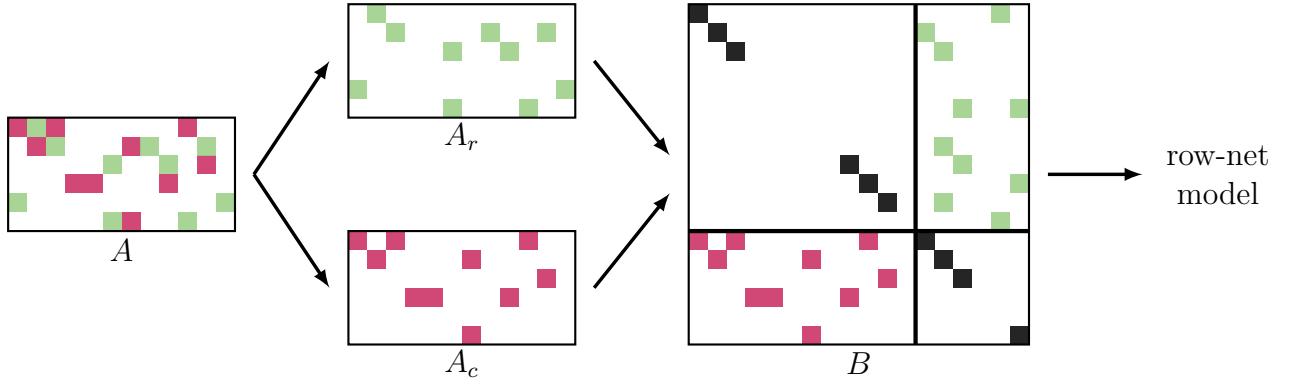


Figure 1.6: Example of the construction of the matrix B from a 6×12 matrix A , for which the sets A_r and A_c were previously established and colored differently. In the resulting matrix, the added diagonal nonzeros are depicted in black. Not all diagonal elements are needed, because of empty rows and columns in A_r and A_c .

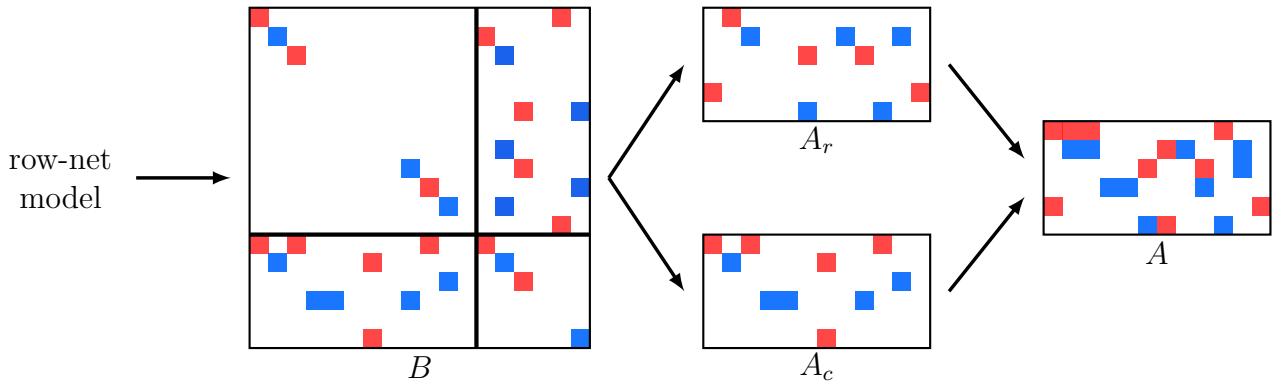


Figure 1.7: Process of obtaining a matrix partitioning starting from a partitioning of the hypergraph following the medium-grain model. In this case $p = 2$.

column of A_c , that belongs to the same row of B as the k th row of A_r , has the purpose of ensuring a correct computation of the communication volume [13, Th. 3.1]. Note that, implementation-wise, there is no need to have the complete diagonal of B : we put a nonzero if and only if the corresponding row of A_r^T and column of A_c are both nonempty.

Experimental results, performed with both the Mondriaan and PaToH packages seem to confirm that this medium-grain method has indeed some advantages compared to the column-net, row-net and fine-grain models, both regarding partitioning time and solution quality. Because of these good results, it is our goal to investigate further the properties of this method, following two possible directions.

First of all, as the outcome of the medium-grain method depends remarkably on the initial split of A into A_r and A_c , it is interesting to investigate the quality of the algorithm originally proposed in [3] to achieve this initial partitioning; secondly, we will try to develop a fully iterative method that employs the medium-grain method, where a full multi-level partitioning is performed at each iteration and computation time is traded for solution quality.

Both these research directions share an important part: we just need to develop efficient methods to compute from the given matrix A the matrices A_r or A_c required for the medium grain model, either from scratch or starting from an already existing partitioning (later in the work we will talk, respectively, about *partition-oblivious* and *partition-aware* algorithms).

To this extent, Chapters 2 and 3 describe several of these different methods, whereas in Chapter 4 we discuss their implementation and the experimental results for the two mentioned research directions.

Chapter 2

Strategies for splitting a matrix A into A_r and A_c

The main goal of this thesis is to find efficient ways of splitting our original matrix A into A_r and A_c , in order to use the medium-grain method.

We are interested in both improving the initial partitioning of A , and a fully iterative method; therefore, we will make the distinction between methods that do not need an initial partitioning (*partition-oblivious* methods), and are therefore suitable for the first case, and methods that do require an initial partitioning (*partition-aware* methods), to be used in a fully iterative scheme. Most of the time the same algorithm can be used for both purposes, albeit with slight modifications. Before we proceed and analyze the details of the examined heuristics, we can make a few observations, to better understand the general principles behind these algorithms.

If we are interested in an initial partitioning into A_r and A_c that will yield a good communication volume, we already have some information about their quality before the actual partitioning is performed. We can indeed compute an upper bound on the communication cost: if a complete row of A is assigned to A_r (or a full column is assigned to A_c), we are sure that those nonzeros will be assigned to the same processor, and we already discussed in Section 1.1 how this results in no communication for that row (or column). This can give us the idea of trying to keep, as much as possible, full rows and columns together, although it is impossible to do it all the time (because a given nonzero cannot be assigned to both A_r and A_c).

If our purpose is to compute A_r and A_c to improve an existing partitioning, we can follow a few principles to guide us in the choice of what information we should keep, and what we should discard for the next iteration. First of all, it makes sense to have confidence in the existing partitioning: if some nonzeros (for example, a full row or column) are assigned to the same processor, it means that at some point in the previous iteration it was decided that it was convenient to put those nonzeros together, and therefore we should have a preference for them to be together also in the new partitioning. However, this must only serve as an indication and not as a rigid rule, leaving some space for new choices to be made, in order to effectively improve the existing partitioning. Furthermore, we should try to keep, as much as possible, rows and columns together, as noted in the previous paragraph.

2.1 Individual assignment of nonzeros

A simple heuristic that can be used to produce A_r and A_c is a simplification of the algorithm proposed by Pelt and Bisseling along with the medium-grain model [13, Alg. 1], taking as a *score* function the length (i.e. the number of nonzeros) of the given row or column.

The main idea is to assign each nonzero a_{ij} to A_r if row i is shorter than column j (so it has a higher

probability of being uncut in a good partitioning), and to A_c otherwise. Ties are broken, similarly as the original algorithm, in a consistent manner: if the matrix is rectangular we give preference to the shorter dimension, otherwise we perform a random choice.

The partition-oblivious version of this heuristic is given in Algorithm 2.1, and it is exactly the same as Algorithm 1 originally proposed. With $nz_r(i)$ we denote the number of nonzeros in row i and with $nz_c(j)$ the number of nonzeros in column j .

Input: sparse matrix A
Output: A_r, A_c

```

if  $m < n$  then
     $w \leftarrow r$ 
else if  $n < m$  then
     $w \leftarrow c$ 
else
     $w \leftarrow$  random value  $c$  or  $r$ 
end if
 $A_r := A_c := \emptyset$ 
for all  $a_{ij} \in A$  do
    if  $nz_r(i) < nz_c(j)$  then
        assign  $a_{ij}$  to  $A_r$ 
    else if  $nz_c(j) < nz_r(i)$  then
        assign  $a_{ij}$  to  $A_c$ 
    else
        assign  $a_{ij}$  to  $A_w$ 
    end if
end for

```

Algorithm 2.1: Partition-oblivious individual assignment of the nonzeros, based on row/column length.

This algorithm can be easily adapted to compute A_r and A_c from a given partitioning of A . Previously we claimed that it is convenient that uncut rows and columns have precedence over cut rows and columns: now, whenever we analyze a nonzero a_{ij} we first look at whether i and j are cut or uncut. If only one of them is cut, we assign the nonzero to the uncut one, otherwise (i.e. both are cut, or both are uncut) we do similarly as before and assign it to the shorter one.

The partition aware variant of this heuristic is given explicitly in Algorithm 2.2.

In Chapter 4, when we are going to perform numerical experiments, we will denote with `po_localview` and `pa_localview`, respectively, the partition-oblivious and the partition-aware variant of this heuristic.

2.2 Assignment of blocks of nonzeros

Instead of assigning nonzeros individually as in Section 2.1, we can take a more coarse-grained approach and try to assign at the same time a block of nonzeros to either A_r or A_c . In particular, we will discuss how to exploit the Separated Block Diagonal (SBD) form of the partitioned matrix A and introduce a further iteration of this concept, discussing the Separated Block Diagonal of order 2 (SBD2) form of the matrix. Moreover, the heuristics described in this section are all partition-aware, and take as input a partitioned matrix.

As throughout this section the permutations of matrices will be fundamental, we adopt a simplified notation: given a vector I with row indices and a vector J with column indices, we denote as $A(I, J)$ the submatrix of A with only the rows in I and only the columns in J (following the order in which they appear in the vectors). With this notation, for example, $A([1, \dots, m], [1 \dots n]) = A$. Furthermore, if I_1 and I_2 are both vectors of indices, with (I_1, I_2) we denote the simple concatenation of these vectors.

Input: partitioned sparse matrix A

Output: A_r, A_c

```

if  $m < n$  then
     $w \leftarrow r$ 
else if  $n < m$  then
     $w \leftarrow c$ 
else
     $w \leftarrow$  random value  $c$  or  $r$ 
end if
 $A_r := A_c := \emptyset$ 
for all  $a_{ij} \in A$  do
    if row  $i$  is uncut and column  $j$  is cut then
        assign  $a_{ij}$  to  $A_r$ 
    else if row  $i$  is cut and column  $j$  is uncut then
        assign  $a_{ij}$  to  $A_c$ 
    else
        if  $nz_r(i) < nz_c(j)$  then
            assign  $a_{ij}$  to  $A_r$ 
        else if  $nz_c(j) < nz_r(i)$  then
            assign  $a_{ij}$  to  $A_c$ 
        else
            assign  $a_{ij}$  to  $A_w$ 
        end if
    end if
end for

```

Algorithm 2.2: Partition-aware individual assignment of the nonzeros, based on row/column length.

2.2.1 Using the Separated Block Diagonal form of A

The SBD form of a bipartitioned matrix [23] is defined as follows: given a matrix A whose nonzeros are either assigned to processor 0 or 1, we compute the vectors R_0 and R_2 of the indices of the rows fully assigned, respectively, to processor 0 and processor 1, and the vector R_1 of the indices of the rows partially assigned to both of the processors; similarly, we compute C_0 , C_2 and C_1 for the columns. Note that, when creating these vectors, their inner ordering is not important; usually, the ascending order is kept.

Then, we obtain the final index vector for the rows as $I = (R_0, R_1, R_2)$ and for the columns as $J = (C_0, C_1, C_2)$. With these quantities, we can finally compute the SBD form of the matrix A as $A(I, J)$.

An example of the procedure for obtaining this form is shown in Figure 2.1.

More explicitly, if we denote as $m_i := |R_i|$, $n_i := |C_i|$, with $i = 0, 1, 2$, the SBD form is the resulting block matrix:

$$\dot{A} := A(I, J) = \begin{bmatrix} \dot{A}_{00} & \dot{A}_{01} & \dot{A}_{02} \\ \dot{A}_{10} & \dot{A}_{11} & \dot{A}_{12} \\ \dot{A}_{21} & \dot{A}_{22} \end{bmatrix}, \quad (2.1)$$

where

- \dot{A}_{00} of size $m_0 \times n_0$, has nonzeros with uncut rows and uncut columns for processor 0;
- \dot{A}_{22} of size $m_2 \times n_2$, has nonzeros with uncut rows and uncut columns for processor 1;
- \dot{A}_{01} of size $m_0 \times n_1$, has nonzeros with uncut rows for processor 0 and cut columns;

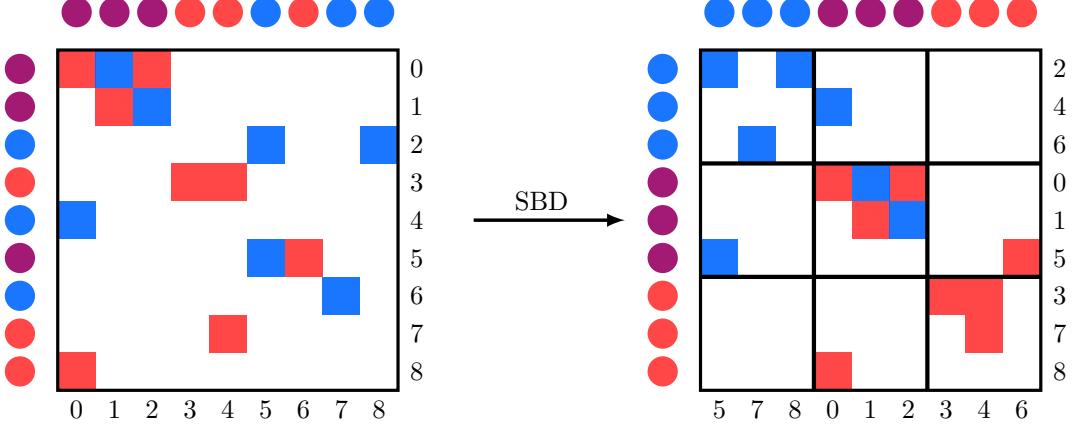


Figure 2.1: Example process to obtain the SBD form of a partitioned matrix. On the left the original matrix is shown, whereas on the right the permuted SBD form. On the top/left sides of the matrices the color of the circle denotes whether that row/column is completely red or blue or it is mixed (purple), whereas on the bottom/right sides the indices of the columns/rows are explicitly given.

- \dot{A}_{21} of size $m_2 \times n_1$, has nonzeros with uncut rows for processor 1 and cut columns;
- \dot{A}_{10} of size $m_1 \times n_0$, has nonzeros with cut rows and uncut columns for processor 0;
- \dot{A}_{12} of size $m_1 \times n_2$, has nonzeros with cut rows and uncut columns for processor 1;
- \dot{A}_{11} of size $m_1 \times n_1$, has nonzeros with cut rows and columns.

Note that the size of each part along with the number of contained nonzeros can greatly vary, also from matrix to matrix: for example, if the sparsity pattern of the matrix allows a “perfect” partitioning such that there is no communication, all blocks are empty except \dot{A}_{00} and \dot{A}_{22} ; conversely, if the matrix has a dense (or complicated) pattern and/or the partitioning is far from optimal, such blocks might be almost empty and \dot{A}_{11} will have the majority of nonzeros. An example of the difference of the block sizes of \dot{A} is shown in Figure 2.2.

By computing the Separated Block Diagonal form of a matrix, we are able to explicitly see the underlying structure of the partitioning of a matrix, and the properties of each block can be used to adapt the assignment of its nonzeros. More specifically, the blocks \dot{A}_{00} and \dot{A}_{22} have nonzeros with uncut rows and columns and therefore are more suited to be assigned together; of course, we still have to decide between A_r and A_c and, as mentioned earlier, it is impossible to do both: it is convenient to base our choice on the sizes of such blocks. For example, if $m_0 < n_0$, in the block \dot{A}_{00} the columns are (on average) sparser than the rows: if we assign the nonzeros of this block to A_c we are, in principle, making sure that more rows/columns will stay uncut.

For the blocks with uncut rows and cut columns (namely, \dot{A}_{01} and \dot{A}_{21}), the choice is easy: we assign them to A_r and keep their rows uncut. Similarly, we assign the nonzeros of \dot{A}_{10} and \dot{A}_{12} to A_c , keeping their columns uncut.

For the middle block \dot{A}_{11} , whose nonzeros have cut rows and cut columns, we cannot exploit any underlying structure: a possible way is to employ one of the other heuristics described in this chapter only considering this submatrix. Our choice is to go with Algorithm 2.1 presented in Section 2.1 (note that we cannot exploit the partition-aware variant of it, because all of the nonzeros in the block considered have cut rows and columns).

The heuristic that employs the SBD structure of a matrix can be explicitly visualized in (2.2).

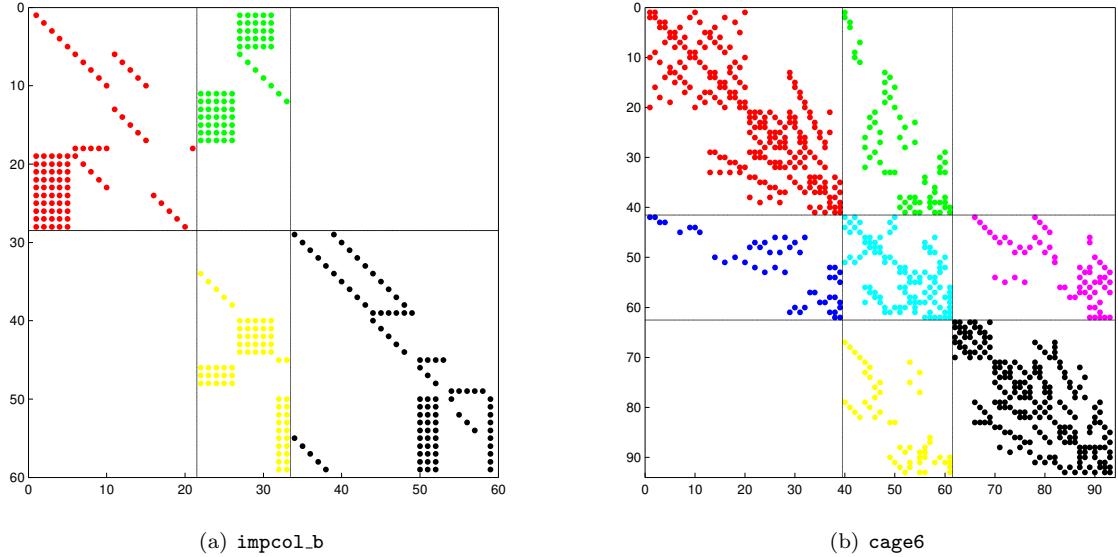


Figure 2.2: Example of SBD forms of partitioning of the matrices `impcol_b` and `cage6` [24]. Each part of \dot{A} has been colored differently. In the first matrix there are no cut rows, which means that $\dot{A}_{10} = \dot{A}_{11} = \dot{A}_{12} = \emptyset$.

$$\text{Assignment of } \dot{A} : \begin{bmatrix} R/C & R \\ C & M & C \\ R & R/C \end{bmatrix} \quad (2.2)$$

In this matrix, whose structure is the same of (2.1), the letter R in a block denotes that we assign that block to A_r , and similarly for C and A_c . Moreover, R/C stands that the choice between A_r and A_c depends on the block size, whereas M indicates that the block is assigned in a mixed manner, according to Algorithm 2.1.

Note that, as mentioned in Chapter 1, the matrix is usually split by means of recursive bipartitionings: it is then sufficient to keep track of the order of these recursions to have an implicit ordering which can be easily used to compute the SBD form of a matrix [23], instead of computing this form from scratch using the algorithm described in [25, Appendix A].

In Chapter 4, we will refer to the algorithm described in this section as `pa_sbdview`.

2.2.2 Using the Separated Block Diagonal form of order 2 of A

The proposed SBD2 form of a partitioned matrix A is an extension of the SBD form: given a partitioned matrix A , we compute the Separate Block Diagonal form of A of order 2 by separating, in \dot{A}_{10} and \dot{A}_{12} the empty and non-empty columns, and in \dot{A}_{01} and \dot{A}_{21} the empty and non-empty rows. Then all the other blocks, except the central one, are permuted and split up accordingly. This procedure is shown in Algorithm 2.3.

The resulting final matrix is a block tridiagonal matrix \ddot{A} :

Input: partitioned matrix A

Output: \ddot{A}

compute \dot{A} as the SBD form of A and obtain also $R_0, R_1, R_2, C_0, C_1, C_2$;
 split R_0 in R_{00} and R_{01} , such that $A(R_{00}, C_1) = \emptyset$;
 split R_2 in R_{20} and R_{21} , such that $A(R_{21}, C_1) = \emptyset$;
 split C_0 in C_{00} and C_{01} , such that $A(R_1, C_{00}) = \emptyset$;
 split C_2 in C_{20} and C_{21} , such that $A(R_1, C_{21}) = \emptyset$;
 $I := (R_{00}, R_{01}, R_1, R_{20}, R_{21})$;
 $J := (C_{00}, C_{01}, C_1, C_{20}, C_{21})$;
 $\ddot{A} := A(I, J)$.

Algorithm 2.3: Algorithm to obtain SBD2 form of a matrix A .

$$\ddot{A} := \begin{bmatrix} \ddot{A}_{00} & \ddot{A}_{01} & & & \\ \ddot{A}_{10} & \ddot{A}_{11} & \ddot{A}_{12} & & \\ & \ddot{A}_{21} & \ddot{A}_{22} & \ddot{A}_{23} & \\ & & \ddot{A}_{32} & \ddot{A}_{33} & \ddot{A}_{34} \\ & & & \ddot{A}_{43} & \ddot{A}_{44} \end{bmatrix}, \quad (2.3)$$

where each submatrix \ddot{A}_{pq} is of size $m_p \times n_q$.

Figure 2.3 shows the process of obtaining this matrix \ddot{A} starting from the SBD matrix \dot{A} obtained in Figure 2.1.

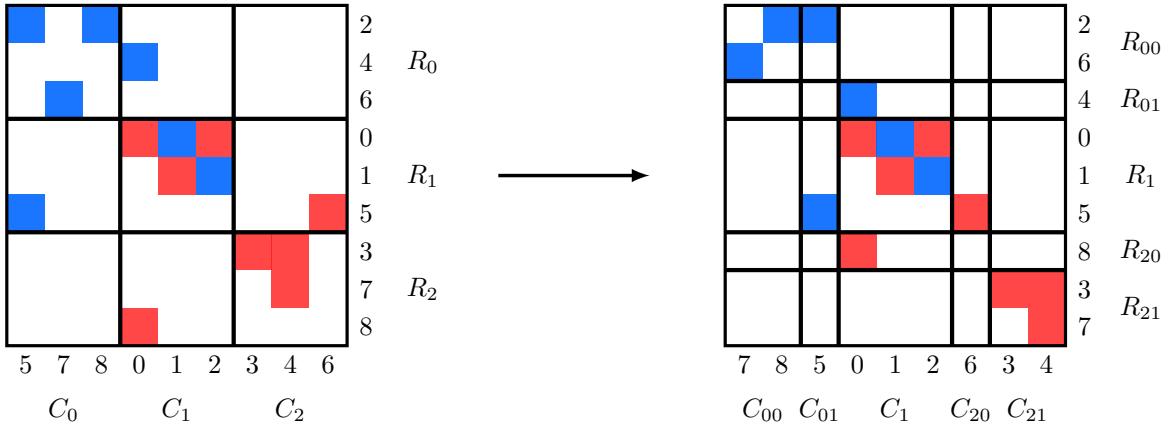


Figure 2.3: SBD2 form obtained starting from the SBD form of Figure 2.1.

To better understand the interesting properties of the newly created parts of the matrix, let us introduce the concept of *neighbor*: given the nonzero a_{ij} we say that a_{kl} is a neighbor if $k = i \vee l = j$; in other words, neighbors of a given nonzero are the ones that lie in the same row or in the same column.

Now, let us consider, for the sake of brevity, just the top-left corner of \ddot{A} : nonzeros in \ddot{A}_{00} are uncut in the rows and columns and whose neighbors are uncut also in the other, non-shared, dimension. Similarly, nonzeros in \ddot{A}_{01} do not have any neighbor (w.r.t. their row) with cut columns but have neighbors (w.r.t their column) with cut rows. And similarly, with the roles of rows and columns reversed, for \ddot{A}_{10} . This exact same reasoning applies also for the bottom-right corner, with the appropriate adaptation of indices.

The size of these parts, and more generally of all of the blocks of \ddot{A} , is again highly dependent on the structure of the matrix, as shown in Figure 2.4.

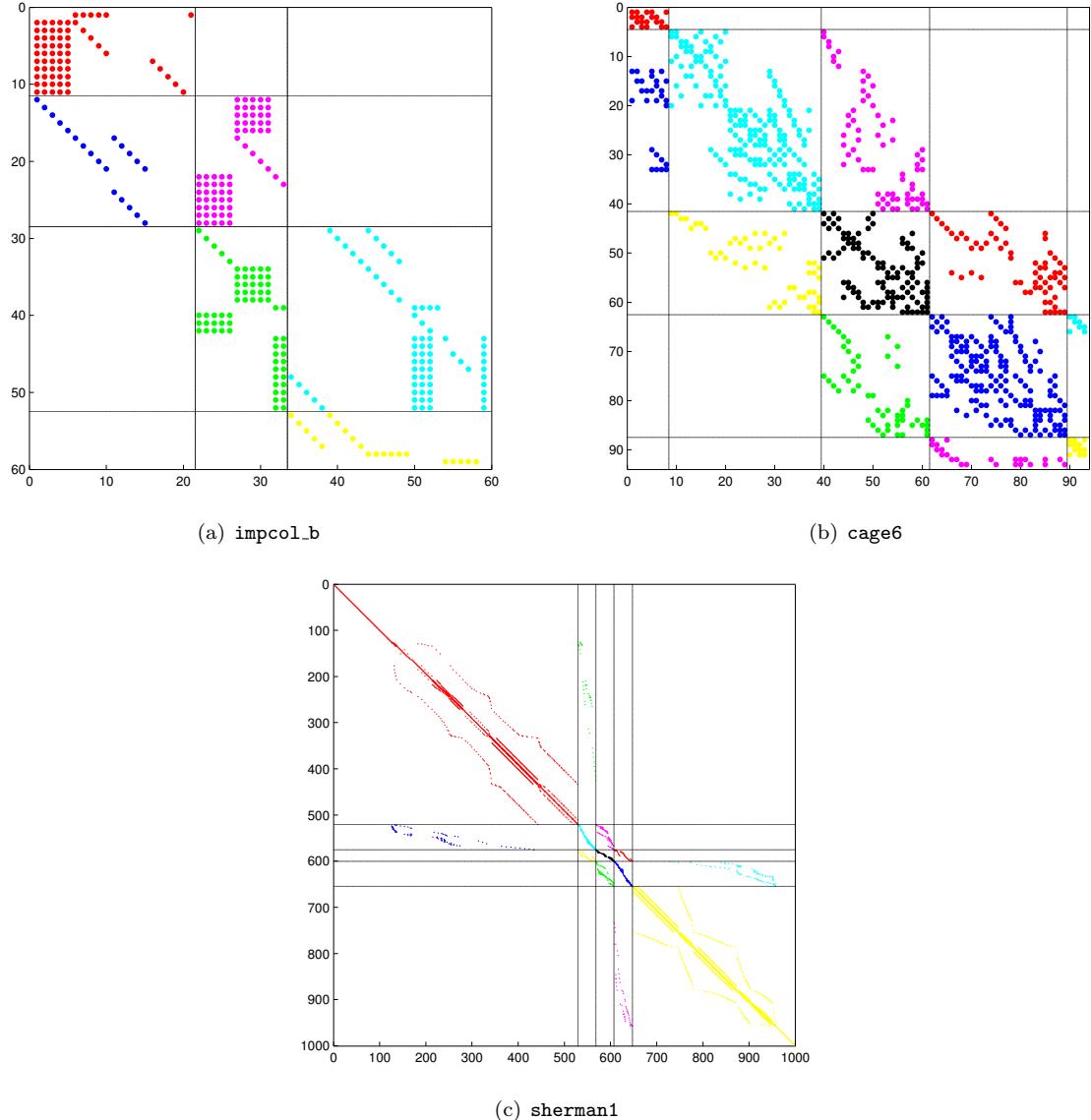


Figure 2.4: Example of SBD2 forms of three different matrices. Similarly as in Figure 2.2, each part of \ddot{A} has been given a color (note that since there are more parts than colors used, some colors are repeated even though the parts are not related in any way). We can see in 2.4(a) that the second and fourth columns are empty, and therefore not shown in the image. We can also see the difference in structure between 2.4(b) and 2.4(c): the former one comes from a DNA electrophoresis problem [24], while the latter is an oil reservoir simulation challenge matrix [26]. We can see that with the `sherman1` matrix, the corner parts are predominant because it is a finite element matrix, with a strongly diagonal pattern: it makes sense that most of these nonzeros are “independent” from each other.

Other than the corner blocks, for which we already argued that the matrix partitioning problem is easy, this structure enables us to assign more specifically nonzeros to either A_r or A_c : it is convenient to assign \ddot{A}_{01} and \ddot{A}_{43} to A_r , as these nonzeros can be fully assigned to one processor without having the columns cut, and similarly we can assign \ddot{A}_{10} and \ddot{A}_{34} to A_c ; for the other blocks, we can repeat the reasoning of the last section.

This heuristic that exploits the SBD2 form of the matrix A is given explicitly as follows:

$$\text{Assignment of } \ddot{A} : \begin{bmatrix} R & R & & \\ C & R/C & R & \\ & C & M & C \\ & & R & R/C & C \\ & & & R & C \end{bmatrix} \quad (2.4)$$

R denotes, as previously, that the block has been assigned to A_r and similarly for C and A_c . R/C denotes that the assignment depends on the block size and M that a mixed assignment is performed with Algorithm 2.1.

Note that, in this case, the SBD2 form has to be computed from scratch from the SBD form, because it uses further information that is not employed during the normal partitioning.

In Chapter 4, we will refer to the algorithm described in this section as `pa_sbd2view`.

2.3 Maximizing empty rows of B

In this section, instead of describing a generating scheme that takes as input the matrix A and produces as output A_r and A_c , we will introduce an improvement scheme, which operates on already existing A_r and A_c and tries to refine them such that the upper bound on the communication volume is lowered.

At the beginning of this chapter, we mentioned how it is convenient to have complete rows assigned to A_r and complete columns assigned to A_c , in order to avoid communication; a good strategy to produce good A_r and A_c , could then be to maximize such complete assignments. The proposed heuristic does essentially this, by trying to swap the assignment of nonzeros from A_r to A_c and viceversa, trying to obtain that complete rows are assigned to A_r and complete columns are assigned to A_c . In order to achieve both of these goals with a unique algorithm, it is convenient to reason in terms of the matrix B as in (1.9). If we maximize the number of empty rows of B , we are effectively emptying rows of A_r^T (i.e. emptying columns of A_r , therefore completely assigning nonzeros in them to A_c) and of A_c , thus complete rows to A_r .

This improvement heuristic falls into the category of *local search* algorithms: we start from a configuration (an assignment of nonzeros to A_r and A_c) and perform a search on the neighborhood, defined as the set of configurations which differ only by the assignment of a single nonzero. By performing this small swap, we can easily fall in a local optimum situation: a few nonzeros (depending on the structure of the matrix) are continuously swapped between A_r and A_c .

We can add a little hill-climbing capability to our heuristic by adding a small buffer: we pre-determine l_{max} , the maximum amount of worsening allowed, and, after this threshold is reached, we start considering only strictly improving solutions. In order to have a meaningful threshold, it might be convenient to have it relative to the number of rows/columns of B , or to its nonzeros. The higher this threshold is, the more capability we have of escaping local optima, but at the cost of slowing down considerably the improvement (even potentially arresting it) of our solution.

For the choice of the neighbor configuration to consider, it is convenient to consider the row of B with a diagonal element (which corresponds to a split row/column of A) with the minimum number of nonzeros: our immediate goal, which in reality spans over a few moves of our local search, is to completely assign the nonzeros of this row of B ; we consider the minimum because each time we swap we might slightly worsen the solution.

A more explicit overview on this local search improvement scheme is described in Algorithm 2.4:

As this is a scheme that relies on existing A_r and A_c and aims at improving them, we still need to choose how to generate these parts in the first place. If we can rely on an existing partitioning, a simple choice could be to take as A_r and A_c the subsets of nonzeros assigned, respectively, to processor 0 and 1; otherwise, if our goal is to produce A_r and A_c for the initial partitioning, the simplest choice is to randomly assign each nonzero to either A_r or A_c . These initial solutions, however, are fast to generate

```

Input:  $A_r$ ,  $A_c$ ,  $l_{max}$ ,  $iter_{max}$ 
Output:  $A'_r$ ,  $A'_c$ 
 $l := 0$ 
Compute initial  $B$  following the medium-grain method
for  $it = 1, \dots, iter_{max}$  do
     $i := \operatorname{argmin}_{k \in \{0, \dots, m+n-1\} \text{ s.t. } B(k,k) \neq 0} nz_r(k)$ 
    for all  $j \neq i$  such that  $B(i,j) \neq 0$  do
        if  $B(j,j) \neq 0$  then
             $B(i,j) = 0$ 
             $B(j,i) = 1$ 
        else
            if  $l < l_{max}$  then
                 $B(i,j) = 0$ 
                 $B(j,i) = 1$ 
                 $l = l + 1$ 
            end if
        end if
    end for
    if  $nz(B(i)) = 1$  then
         $B(i,i) = 0$ 
         $l = l - 1$ 
    end if
end for
 $A'_r := B([0, \dots, n-1], [n, \dots, m+n-1])^T$ 
 $A'_c := B([n, \dots, m+n-1], [0, \dots, n-1])$ 

```

Algorithm 2.4: Local search refinement of A_r and A_c

but not particularly efficient, and are therefore meaningful only if our improvement scheme is fast enough; otherwise, we can always rely on one of the other heuristics described in this chapter.

2.4 Partial assignment of rows and columns

In Section 2.1 we discussed how to assign each nonzero independently, whereas in Section 2.2 we examined the possibility of exploiting a little the structure of the matrix, in order to assign more nonzeros at once. Keeping this direction, there is some other structure of A that can lead to a better assignment: partial assignment of rows and columns.

The main idea behind this heuristic is that, every time we assign a nonzero to A_r , we know that it is convenient that also all the other nonzeros in the same row are assigned to it; conversely, if a nonzero is assigned to A_c , all the nonzeros in its column should stick with it. Therefore, we ideally want to keep together rows and columns as much as possible; but, as already discussed, there is always the problem that a nonzero cannot be assigned to both A_r and A_c , so that we can only reason in term of partial assignment of the row/column.

Throughout this section we will stop distinguishing between rows and columns of a matrix and reason in term of **indices** in the set $\{0, \dots, m+n-1\}$: following the natural ordering, the m rows are mapped to $0, \dots, m-1$ and the n columns to $m, \dots, m+n-1$.

This simplification of terms is due to the fact that the core of this heuristic lies in the computation of a **priority vector** v , which is none other than a permutation of the indices $0, \dots, m+n-1$, where they appear in order of decreasing priority: in this sense, the priority is to be intended as the probability of the nonzeros of that index to be together in a good partitioning.

The assignment of nonzeros is done by “painting” them with an imaginary color, which corresponds

either to A_r or A_c : we iterate through our priority vector backwards (i.e. starting from the index with the lowest priority) and assign all of its nonzeros to go together: if the index corresponds to a row, then we assign all of its nonzeros to A_r , otherwise we assign them to A_c . Because each nonzero has both a row and a column, it is represented twice in our priority vector; the second time it is considered, we re-assign it by “painting it over” (hence the name of the algorithm).

A more explicit formulation of this procedure is given in Algorithm 2.5.

```
Input: Priority vector  $v$ , matrix  $A$ 
Output:  $A_r, A_c$ 
 $A_r := A_c := \emptyset$ 
for  $i = m + n - 1, \dots, 0$  do
  if  $v_i < m$  then
    Add the nonzeros of row  $i$  to  $A_r$ 
  else
    Add the nonzeros of column  $i - m$  to  $A_c$ 
  end if
end for
```

Algorithm 2.5: Overpainting algorithm

It is possible also to give an alternative formulation for this algorithm in which we iterate forward through v , as described in Algorithm 2.6.

```
Input: Priority vector  $v$ , matrix  $A$ 
Output:  $A_r, A_c$ 
 $A_r := A_c := \emptyset$ 
for  $i = 0, \dots, m + n - 1$  do
  if  $v_i < m$  then
    Add the unmarked nonzeros of row  $i$  of  $A$  to  $A_r$ 
  else
    Add the unmarked nonzeros of column  $i - m$  of  $A$  to  $A_c$ 
  end if
  Mark nonzeros of index  $i$  as “evaluated”
end for
```

Algorithm 2.6: Alternative formulation of Algorithm 2.5.

In this formulation, every assignment to A_r and A_c is final, but with the added complexity of checking which nonzeros of the considered index are still to be assigned, and only working with them.

Lastly, another different formulation is possible: we consider individually each nonzero a_{ij} and see whether in v $i < j$ (where the $<$ symbol is to be intended as “ i precedes j ” and not as the comparison of the values) or the other way around; in the first case, the row has more priority and we assign a_{ij} to A_r , otherwise we assign it to A_c . Note that, since we have to perform N lookups on the vector v , this is a more expensive formulation of the same algorithm.

An important point to observe is that this overpainting algorithm is completely deterministic: A_r and A_c are uniquely determined by the ordering of the indices in v . Therefore, the heuristic part of this algorithm lies entirely in the choice of this priority vector, and, for this reason, we will focus on it in the next subsection.

2.4.1 Computation of the priority vector v

Because, with the overpainting algorithm, the quality of A_r and A_c depends entirely on the choice of v , it is important to take a structured approach and explore a wide variety of possibilities for this priority vector.

In this section, we proceed and define several *generating schemes*, and their input determines whether the overpainting algorithm is used to obtain a better initial partitioning or a fully iterative scheme. In general, we try to come up with schemes that can be used for either purpose, with slight modifications.

Each one of the generating schemes can be summarized in three main steps:

1. usage of previous partitioning;
2. sorting;
3. internal ordering of indices.

Now, we give a more detailed explanation of each of those steps.

- **usage of previous partitioning:** if we are considering *partition-aware generating schemes*, we separate the set of uncut indices (i.e. indices which correspond to uncut rows or columns) from the cut indices. We consider the simple concatenation of uncut indices and cut indices, in this order, and the next steps are performed on each of these parts. If, instead, we are considering a *partition-oblivious scheme*, the subsequent operations are performed on the set $\{0, \dots, m+n-1\}$.
- **sorting:** we can either keep the set from the previous step untouched (therefore preserving the natural order of indices) or perform a sorting with respect to the number of nonzeros. The sorting is done in ascending order, as a short row/column is more likely to fit completely in a good partitioning because it does not yield many cut columns/rows.
- **internal ordering:** as the last step, we want to finalize our vector v by deciding more precisely the position of each index. The strategies considered, which often depend internally on an additional parameter, are the following:
 - **concatenation:** we put either all the rows before all the columns, or all the columns before all the rows;
 - **mixing:** we can mix rows and columns in two main ways: alternation and spread. Suppose there are twice as many columns as rows: in the first case we get

$$(c, r, c, r, c, r, \dots, c, r, c, c, c, \dots, c, c, c),$$

whereas with the second one we get

$$(c, c, r, c, c, r, \dots, c, c, r),$$

where with c we denote a generic column and with r a generic row. To obtain a more even distribution, we always start with the greater dimension.

- **random** (only in case of no sorting): we randomize the ordering of the indices;
- **simple** (only in case of sorting): we let the sorting decide completely the ordering, and the vector is left as-is.

As the complete description of a generating scheme is somewhat lengthy, we use a simplified notation and adopt the following abbreviations:

- P0: partition-oblivious
- PA: partition-aware

- **sorted** and **unsorted**: sorting w.r.t. the number of nonzeros is performed or not
- **w** and **nw**: all the indices with only 1 nonzero are moved to the back or not
- **simple**: the sorted vector is left as-is
- **concat**: rows and columns are concatenated
- **row**: the concatenation is done rows-columns
- **col**: the concatenation is done columns-rows
- **mix**: mixing of the rows and columns is enforced
- **alt**: rows and columns are alternated
- **spr**: rows and columns are spread
- **random**: the order of the indices is randomized

With this notation, the name `po_sorted_nw_mix_spr` stands for “partition-oblivious generating scheme, with indices sorted by number of nonzeros, without moving the indices with 1 nonzero to the back, with forced mixing of rows and columns, in a spread fashion”. This is just one of the many possibilities, which are convenient to visualize using a directed graph, as shown in Figure 2.5; a generating scheme is simply a path from START to END.

Other than this family of heuristics, we can also formulate the problem of partial assignment of rows/columns, always following this framework, in another more mathematical way, to which we dedicate Chapter 3.

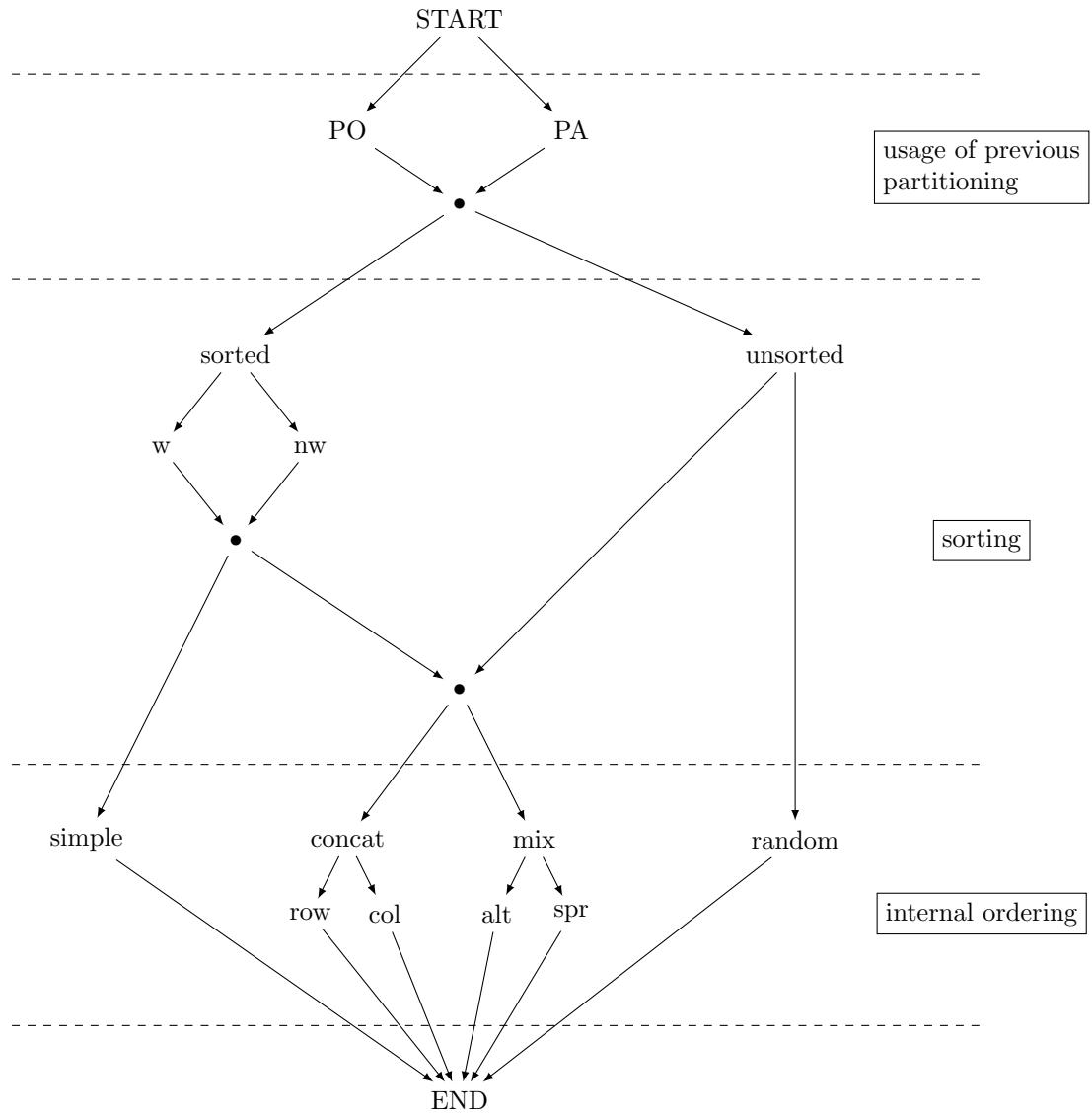


Figure 2.5: Directed graph that represents the family of heuristics used (any path from START to END). Dummy nodes (the ones without any label) were added in order to reduce the number of edges and ease legibility.

Chapter 3

Maximum independent set formulation of the partial row/column assignment problem

With the framework introduced in Section 2.4, we basically translated the problem of the assignment of nonzeros to A_r and A_c (which is already another formulation of the matrix partitioning problem with the medium grain model) to the problem of an efficient computation of a permutation of the indices $\{0, \dots, m+n-1\}$. In this chapter, we will propose a method for this vector computation problem which relies on concepts of the field of graph theory.

The main idea is somewhat similar to the principle that led us to the development of the Separated Block Diagonal form of order 2 in Section 2.2.2. In that particular form of a partitioned matrix, the blocks \ddot{A}_{00} and \ddot{A}_{44} are interesting, as they contain “independent” nonzeros. More specifically, those rows and columns are fully assigned to a processor, and whose nonzeros do not have any neighbor (a nonzero in the same row or column) which has a cut column/row. The analogue of this concept of independence, is now to be defined carefully: we want to find a subset of the indices $\{0, \dots, m+n-1\}$ which does not cause any communication, whenever we fully assign its rows to A_r and its columns to A_c . With this definition, our goal is clear: we want to assign as many nonzeros as possible in this way, obtaining a low upper bound on the communication volume, which can be computed during the creation of A_r and A_c .

To do so, we can employ a well studied object in graph theory: the **maximum independent set**. However, this requires a correct translation of our sparse matrix into a graph, described in Section 3.1. In Section 3.2, we delve a little more into the graph theory required and describe the actual algorithm used to compute a maximum independent set in such a graph. In Section 3.3, finally, we give a few different possibilities for computing the priority vector v using the concepts and algorithms just introduced.

3.1 Graph construction

We need to construct the graph correctly from our sparse matrix, in order to retrieve the desired information. In our case, we can simply consider the graph whose adjacency matrix is none other than the sparsity pattern of our matrix A . This exact same formulation has already been studied, for example, by Hendrickson and Kolda [27], who used their *bipartite graph model* to discuss different algorithms for bipartite graph partitioning.

More explicitly, in this graph formulation, rows and columns are vertices, and we have the edge (i, j) if $a_{ij} \neq 0$. It is fairly clear that the resulting graph is bipartite, because an edge connects only a row with a column.

An example of such translation from matrix to graph is shown in Figure 3.1, where we start from the

matrix given in Figure 1.1.

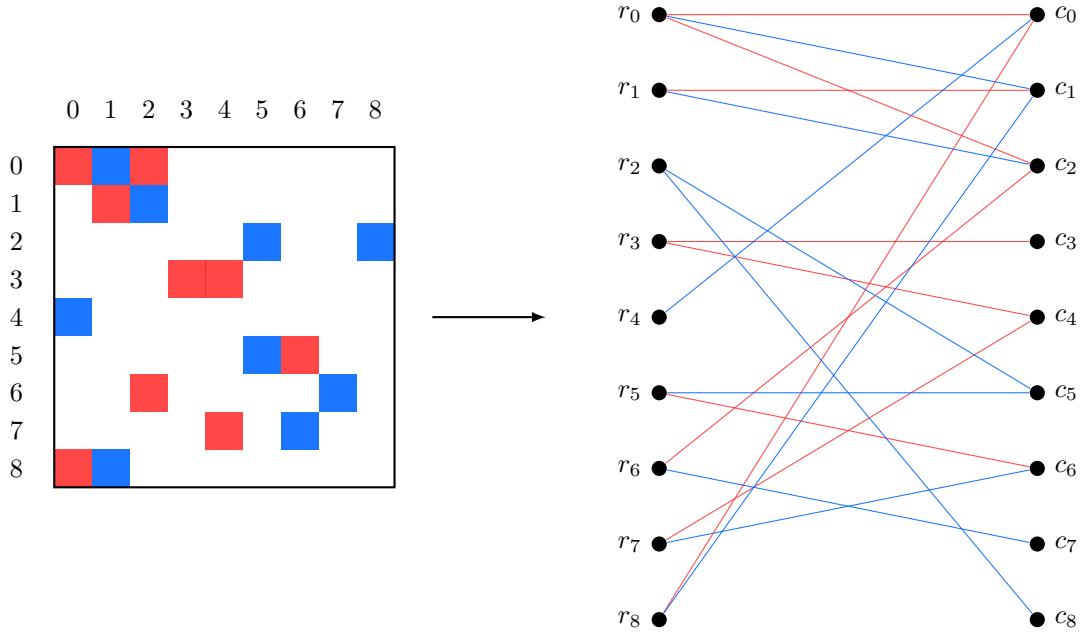


Figure 3.1: Graph constructed using the sparsity pattern of the matrix of Figure 1.1 as adjacency matrix (rows and columns are vertices, nonzeros are edges). The edge color has been kept the same as the corresponding nonzero, but only to facilitate the understanding. The fact that the matrix is partitioned does not play any role in the resulting graph. In the bipartite graph, with r_i we denote row i , whereas with c_j we denote column j .

3.2 The maximum independent set and its computation

In this section, we will give an overview of the maximum independent set problem, discuss its complexity and the relation with other famous problems in graph theory, and, lastly, give an efficient algorithm that can be used with a bipartite graph.

3.2.1 Maximum independent set

The concepts of *independent set* and *vertex cover* are closely related [28]: let $G = (V, E)$ be an undirected graph.

Definition 3.1 (Independent set). *An independent set is a subset $V' \subseteq V$ such that $\forall u, v \in V'$, $(u, v) \notin E$. A maximum independent set is an independent set of G with maximum cardinality.*

Definition 3.2 (Vertex cover). *A vertex cover is a subset $V' \subseteq V$ such that $\forall (u, v) \in E$ we have $u \in V' \vee v \in V'$, i.e. at least one of the endpoints of any edge is in the cover. A minimum vertex cover is a vertex cover of G with minimum cardinality.*

A graphical depiction of two independent sets for an example graph is shown in Figure 3.2.

The following lemma explicitly gives us the relation between a vertex cover and an independent set.

Lemma 3.1. *Given a graph G , V' is a vertex cover set if and only if $V \setminus V'$ is an independent set.*

Proof. Let V' be a vertex cover, i.e. $\forall (u, v) \in E$, $u \in V'$ or $v \in V'$. This is equivalent to say that $\forall u, v \in V \setminus V'$, $(u, v) \notin E$, which is the definition of independent set. \diamond

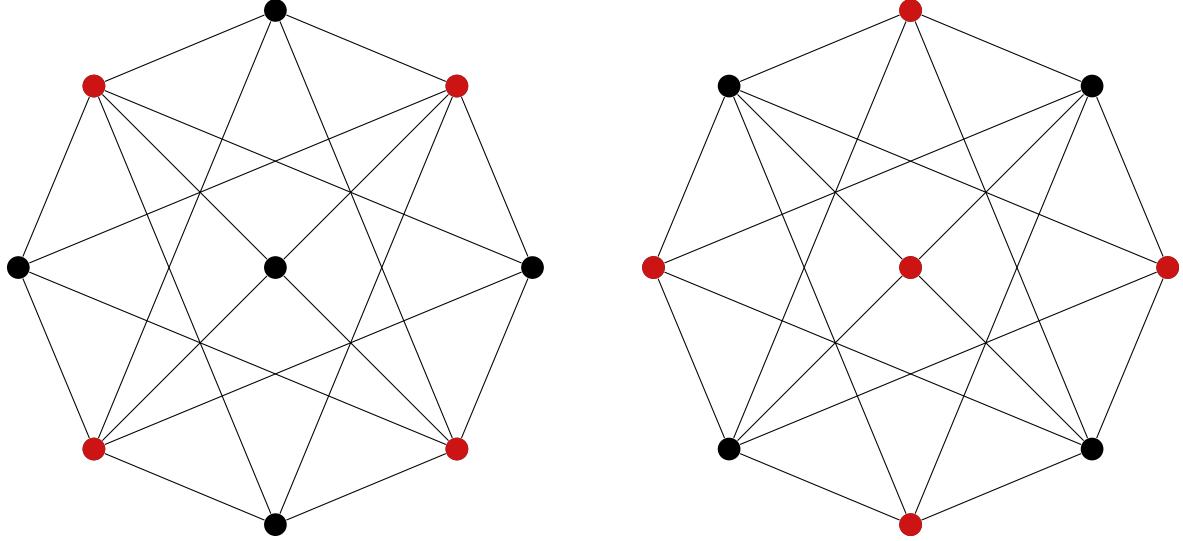


Figure 3.2: Two different independent sets (red vertices) on an example graph. The two independent sets have different cardinality, and the one on the right is a maximum independent set.

As the decision variant of the problem of finding a minimum vertex cover is NP-complete [28, Theorem 3.3], it follows from this lemma that also finding a maximum independent set in a graph is NP-complete; the main consequence of this result is that we cannot solve this problem directly for a generic graph, as it would be as hard as our original matrix partitioning problem. Luckily, we are dealing with a particular kind of graph, a bipartite graph, which simplifies greatly the computations of a maximum independent set.

Before exploiting the bipartiteness of our graph, we need to make an additional observation: Lemma 3.1 states that, in a generic graph, the vertex cover problem and independent set cover are complementary. Therefore, computing a maximum independent set is equivalent to computing a minimum vertex cover. This equivalence is particularly useful in our case, because another mathematical object can be related to the minimum vertex cover as well: the maximum matching.

Definition 3.3 (Matching). *Let $G = (V, E)$ be a graph. A matching $M \subseteq E$ is a set of edges such that at most one edge from M is incident to each vertex $v \in V$. We say that a vertex $v \in V$ is matched by M if an edge in M is incident to v . A maximum matching is a matching of maximum cardinality.*

In particular, because we are in a bipartite graph, we can employ Kőnig's Theorem [29]:

Theorem 3.1 (Kőnig). *In a bipartite graph, the size of a maximum matching is equal to the size of a minimum vertex cover.*

We will algorithmically prove the theorem, showing that from a maximum matching we can obtain a minimum vertex cover and their size is equal. First, however, we need two more definitions that are useful when dealing with (maximum) matchings.

Definition 3.4 (Simple path). *Let $G = (V, E)$ be a graph. A path $P = (v_1, \dots, v_k)$ is said to be simple if $v_i \neq v_j$, $\forall i \neq j$, i.e. all the vertices are distinct and there are no self-edges or sub-cycles.*

Definition 3.5 (Augmenting path). *Let M be a matching on the graph $G = (V, E)$. The simple path P is said to be augmenting if it starts and ends on unmatched (or exposed) vertices, and its edges alternate between $E \setminus M$ and M , in this order.*

It is easy to see that, if we have a matching M and an augmenting path P , if P contains k edges in M , then it has exactly $k + 1$ edges in $E \setminus M$, and, if the graph is bipartite, the two endpoints of P belong to the two different sets of vertices. Moreover, $M \oplus P$ is a matching of size $|M| + 1$, where $M \oplus P := (M \setminus P) \cup (P \setminus M)$ denotes the symmetric difference between M and P .

Now, in Algorithm 3.1, we give a scheme that constructs a bipartite directed graph, starting from a bipartite undirected graph $G = (L \cup R, E)$ and a matching M .

Input: Bipartite undirected graph $G = (L \cup R, E)$, matching M .

Output: Bipartite directed graph $G' = (L \cup R, D)$.

```

 $D \leftarrow \emptyset$ 
for  $(i, j) \in E$  such that  $i \in L$  and  $j \in R$  do
    if  $(i, j) \in M$  then
         $D \leftarrow D \cup (j, i)$ 
    else
         $D \leftarrow D \cup (i, j)$ 
    end if
end for
```

Algorithm 3.1: Construction of a bipartite directed graph starting from an undirected bipartite graph and a matching.

In other words, the edges in E are given a direction: the ones in the matching go from R to L , and the others from L to R .

In Algorithm 3.2, we give an explicit scheme to compute a minimum vertex cover starting from a maximum matching.

Input: Bipartite graph $G = (L \cup R, E)$, maximum matching M .

Output: Minimum vertex cover C .

Construct the modified graph $G' = (L \cup R, D)$ as in Algorithm 3.1

```

 $T \leftarrow \emptyset$ 
for all  $v \in L$  such that  $v$  is not matched do
    Add  $v$  to  $T$ 
    for all  $u \in L \cup R$  reachable from  $v$  using the directed edges in  $D$  do
        Add  $u$  to  $T$ 
    end for
end for
 $C \leftarrow (L \setminus T) \cup (R \cap T)$ 
```

Algorithm 3.2: Construction of the minimum vertex cover in a bipartite graph, starting from the maximum matching.

In Figure 3.3, we can visualize an example of Algorithm 3.2, starting from the graph of Figure 3.1, with a maximum matching. In the final image, the relationships between maximum independent set, minimum vertex cover and maximum matching can be easily recognized.

With the following lemma, we will prove the correctness of Algorithm 3.2 and König's Theorem.

Lemma 3.2. *Let $G = (L \cup R, E)$ be a bipartite graph and let M be a maximum matching. The subset of vertices C obtained following Algorithm 3.2 is the minimum vertex cover. In addition, we have that $|C| = |M|$.*

Proof. First of all, we will prove that C is a vertex cover: assume it is not, which means that there is an edge $e = (i, j) \in E$ with $i, j \notin C$ ($i \in L$, $j \in R$), which implies that $i \in L \cap T$ and $j \in R \setminus T$, because the graph is bipartite.

We now have two possibilities: either $e \notin M$ or $e \in M$. In the first case, because $i \in T$ and e can be traversed (it goes from L to R in the directed graph), j can be reached and thus $j \in T$, a contradiction. In the second case, we have that i is matched and belongs to T : this means that it was reached by traversing the matched edge, which implies that $j \in T$, again a contradiction. Since in both cases we get to a contradiction, C is indeed a vertex cover of G .

Now, we will prove that $|C| \leq |M|$, by showing that every vertex in C is matched. It is clear that every vertex in $L \setminus T$ is matched, by definition of T . Suppose there is an unmatched vertex $v \in R \cap T$: since

$v \in T$, it means that it was reached from an unmatched vertex u in L , thus the path from u to v is augmenting, which would imply that the matching M is not of maximum cardinality, a contradiction. Moreover, note that there are no edges in the matching between $L \setminus T$ and $R \cap T$ (otherwise, as shown previously, the endpoint in L would also be in T).

Therefore we have that all the vertices in C are matched and the edges of the matching are distinct, which implies that $|C| \leq |M|$. Furthermore, for any matching M' and vertex cover C' it is true that $|M'| \leq |C'|$ as there is at least one endpoint in C' for every edge in M' . So we have that $|C| = |M|$, proving Kőnig's Theorem, along the way.

Lastly, these inequality also imply that C is the vertex cover of minimum cardinality: assume it is not, i.e. we have that C' is a vertex cover with $|C'| < |C|$. Now, if we consider the maximum matching M , the inequalities give us that $|M| \leq |C'| < |C| \leq |M|$, a contradiction. C is then the minimum vertex cover.

◊

We have shown that there are close relationships between the maximum matching, minimum vertex cover and maximum independent set. Now, the problem is shifted toward finding an efficient way of computing the maximum matching. In the next section we will describe the Hopcroft-Karp algorithm, which is a simple extension of Algorithm 3.2.

3.2.2 The Hopcroft-Karp algorithm for bipartite matching

The Hopcroft-Karp algorithm [30], devised in 1973, is an efficient scheme for finding a maximum independent set on bipartite graphs, with a running time of $\mathcal{O}(|E|\sqrt{|V|})$. This is a considerable improvement over the famous Ford-Fulkerson algorithm of 1956, which, for bipartite graphs, has a running time of $\mathcal{O}(|V||E|)$. We can compare these two algorithms, even though the latter is technically meant for maximum flow problems, because a bipartite graph can be modified in such a way that a maximum flow corresponds to a maximum matching in the original graph.

Both algorithms rely on the concept of augmenting paths, introduced in the previous section in Definition 3.5.

The main idea of the Hopcroft-Karp algorithm is to find these augmenting paths to progressively increase the size of the matching, as outlined in Algorithm 3.3. The fact that in the main loop we augment the matching over several augmenting paths simultaneously gives us the $\sqrt{|V|}$ factor in the running time, instead of a simple $|V|$.

Input: Bipartite graph $G = (L \cup R, E)$
Output: Maximum matching M

```

 $M \leftarrow \emptyset$ 
repeat
     $l_M \leftarrow$  length of the shortest augmenting path, using the matching  $M$ 
     $P \leftarrow \{P_1, \dots, P_k\}$ , a maximal set of vertex-disjoint shortest augmenting paths of length  $l_M$ 
     $M \leftarrow M \oplus (P_1 \cup \dots \cup P_k)$ 
until  $P = \emptyset$ 

```

Algorithm 3.3: Basic outline of the Hopcroft-Karp algorithm

The core of this algorithm is substantially Algorithm 3.2: instead of starting from a maximum matching M to compute the set T (from which follows immediately a minimum vertex cover C), we construct the matching and the set T progressively, as follows:

1. we construct the directed graph as in Algorithm 3.1;
2. we perform a breadth-first search (following the directed edges) starting from the unmatched vertices in L , which terminates when unmatched vertices in R are reached. l_M is the length of these

shortest augmenting paths;

3. the maximal set of vertex-disjoint shortest augmenting paths is computed: we start from an unmatched vertex in R reached in the previous step and perform a depth-first search. Whenever we reach an unmatched vertex in L it means that we found an augmenting path P , because we were following the directed edges constructed in the first step. We then add this path to P and resume with the next depth-first search.

If we use the Hopcroft-Karp algorithm in our sparse graph constructed as in Section 3.1, the running time can be even considerably better than the theoretical one: if there are no particularly dense rows and columns, the graph is far from being strongly connected, as each vertex in the graph has just a handful of edges, resulting in fast search phases.

3.3 Computation of the priority vector v with the maximum independent set

After having translated our matrix into a graph as in Section 3.1 and having computed the maximum independent set as described in Section 3.2, we still have to compute our priority vector v , to be used in the same framework of Section 2.4. Similarly as done for all the methods described in Chapter 2, we will distinguish between partition-oblivious heuristics and partition-aware ones.

Let $I \subseteq \{0, \dots, m+n-1\}$ be a set of indices. Instead of computing the graph starting from the full matrix A , we do it from the submatrix $A(I)$ (i.e. only taking rows and columns in I); next, we compute the maximum independent set on the resulting graph using the Hopcroft-Karp algorithm: if we denote by S_I the indices that correspond to this maximum independent set, we always give to this set a high priority, putting it before the remaining indices of $I \setminus S_I$.

With this in mind, the partition-oblivious version is quite straightforward: we take as $I = \{0, \dots, m+n-1\}$, and simply compute

$$v := (S_I, I \setminus S_I).$$

Now, for a partitioned matrix, let U denote the set of uncut indices, and C the set of cut indices. For the partition-aware version of this heuristic we have the following possibilities:

1. we compute S_U and have

$$v := (S_U, U \setminus S_U, C);$$

2. we compute S_U, S_C and have

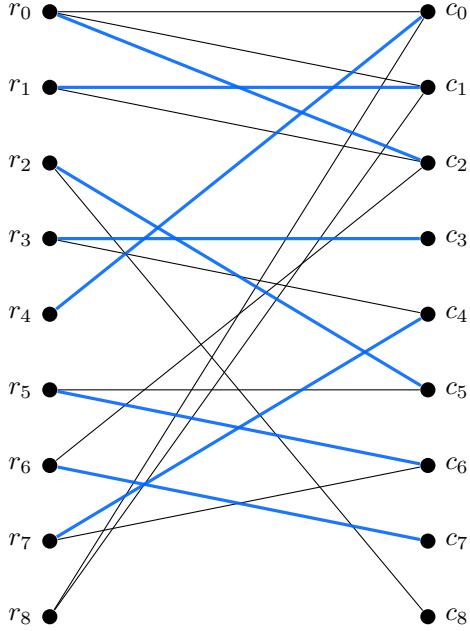
$$v := (S_U, U \setminus S_U, S_C, C \setminus S_C);$$

3. we compute S_U , then we define $U' := U \setminus S_U$ and compute $S_{C \cup U'}$, having

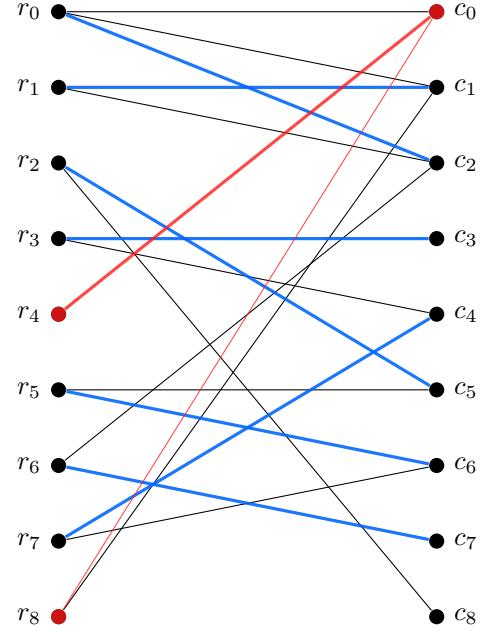
$$v := (S_U, S_{C \cup U'}, (C \cup U') \setminus S_{C \cup U'}).$$

Note that, by construction, we do not expect these three strategies to be radically different in practice: if at the previous iteration the partitioning was done well, U will be quite big, resulting in similar priority vectors.

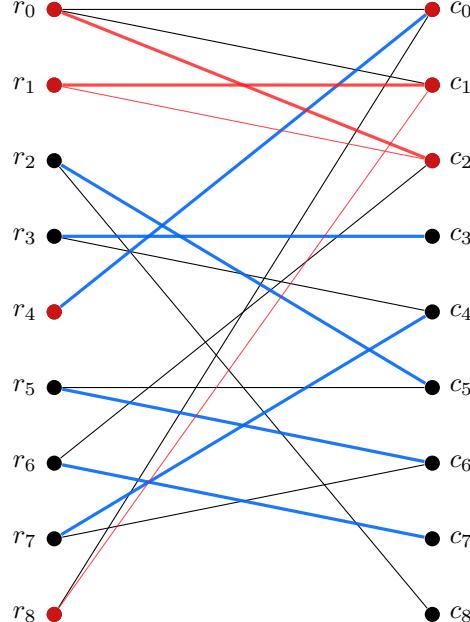
In Chapter 4 we will refer to these three different heuristics, respectively, as `po_is_1`, `po_is_2` and `po_is_3`.



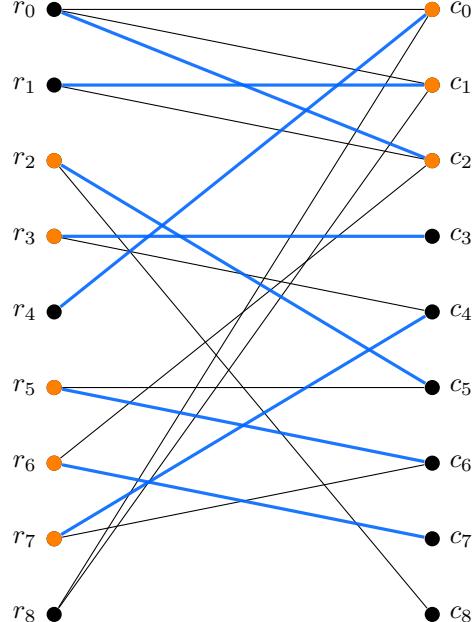
(a) Graph at the beginning of the algorithm. The edges in blue belong to the maximum matching M . $T = \emptyset$.



(b) We start with the unmatched vertex in \$r_8\$; we traverse the unmatched edge \$(r_8, c_0)\$ and the matched edge \$(c_0, r_4)\$. $T = \{r_8, c_0, r_4\}$.



(c) We traverse the unmatched edge \$(r_8, c_1)\$, then the matched edge \$(c_1, r_1)\$, then the unmatched \$(r_1, c_2)\$ and lastly the matched \$(c_2, r_0)\$. $T = \{r_8, c_0, r_4, c_1, r_1, c_2, r_0\}$.



(d) We take as $C = (L \setminus T) \cup (R \cap T)$ (depicted in orange). We can see from the final graph that it is indeed a minimum vertex cover.

Figure 3.3: Example of the actions performed in Algorithm 3.2 for the graph of Figure 3.1. The red vertices are in the set T . In the last image the orange vertices belong to a minimum vertex cover. Lemma 3.1 can quickly be checked, as the black vertices are indeed a maximum independent set.

Chapter 4

Implementation and experimental results

In Chapter 2 and 3 we discussed different heuristics aimed at solving the matrix partitioning problem; whether we try to improve the initial partitioning or perform a fully iterative procedure, we need to translate those ideas into practice, devising an efficient implementation.

In Section 1.3, we mentioned existing software partitioners: Mondriaan [3] is the package of our choice and we defer to it the actual computations of the partitionings, limiting ourselves to create the matrix B of the medium-grain model as in (1.9). The actual algorithm used to construct this matrix is given explicitly in Algorithm 4.1. Note how we consider only the sparsity patterns of A , neglecting completely the values of the nonzeros.

Input: A_r, A_c
Output: B
$B \leftarrow \emptyset$ for all $a_{ij} \in A_c$ do ▷ The part relative to A_c $b_{i+n,j} = 1$ end for for all $a_{ij} \in A_r$ do ▷ The part relative to A_r $b_{j,i+n} = 1$ end for for $j = 1, \dots, n$ do ▷ Dummy nonzeros for cut columns if $\exists i$ s.t. $a_{ij} \in A_r$ and $\exists i'$ s.t. $a_{i'j} \in A_c$ then $b_{j,j} = 1$ end if end for for $i = 1, \dots, m$ do ▷ Dummy nonzeros for cut rows if $\exists j$ s.t. $a_{ij} \in A_r$ and $\exists j'$ s.t. $a_{ij'} \in A_c$ then $b_{n+i,n+i} = 1$ end if end for

Algorithm 4.1: Construction of B following the medium-grain model.

Now, having discussed the means of obtaining A_r and A_c and the matrix B , we can outline the general framework used to test the effectiveness of the proposed heuristics. The framework is given explicitly in Algorithm 4.2, and it takes as a parameter the maximum number of iterations allowed, $iter_{max}$.

In Chapter 2 and 3, we distinguished between partition-oblivious and partition-aware methods, and Algorithm 4.2 is suitable for both types of heuristics: even though the framework is naturally suited for developing a fully iterative scheme, if we desire a better initial partitioning for the medium-grain

Input: Sparse matrix A

Output: Partitioning for the matrix A

Partition A with Mondriaan using the default options and the medium-grain method

for $i = 1, \dots, iter_{max}$ **do**

 Use any of the heuristics described previously to compute A_r and A_c

 construct B , using Algorithm 4.1, from A_r and A_c

 Partition B with Mondriaan using the default options and the row-net model

 Re-construct A with the new partitioning

end for

Algorithm 4.2: General framework for the testing of our heuristics

method, we can simply neglect the partitioning done in the first step. This is precisely the scope of a partition-oblivious heuristic.

Regarding the actual implementation, we can see from Algorithm 4.2 that Mondriaan is used to perform the actual partitioning and this is the ideal case for its use as a software library. As a consequence, we used C as the main implementation language, even though MATLAB was used for faster prototyping: the flexibility added by managing objects at runtime is ideal when designing algorithms. In order to have C code and MATLAB code interact in the correct way, we took advantage of MEX files [31]. In general, unless preliminary tests showed that the considered heuristic had a remarkably bad quality, we translated back most of the programs to the C language, in order to remove the MEX layer of complexity and get a more efficient implementation. For the Hopcroft-Karp algorithm described in Chapter 3, we used an implementation [32] written in the Python programming language, which computes directly the matching on a bipartite graph and the maximum independent set.

In the following tests, the parameter $iter_{max}$ has been set to 1, which means that we are only performing one iteration of our heuristic. The reason for this choice will become clear in Section 4.3; for now, we will just take advantage of the fact that performing only one iteration is relatively fast and we can quickly get an idea about the behavior of the heuristics.

In order to perform effective numerical experiments, we need to have a consistent way of testing. First of all, since we seek heuristics suitable for many matrices, it makes sense to have several matrices to test for, as described in Section 4.1.

Secondly, since randomness is involved in the partitioner itself and, to a different extent, in some of the heuristics, we need to take several measurements and compute an average. We generate 20 independent initial partitionings and, for each of them, 5 independent iterations of the heuristic are performed and averaged. The values used for our measurements and considerations are the average of the 20 initial partitionings and the average of the 20 averages of final partitionings.

Lastly, in order to have a meaningful result that can help us understand whether the given heuristic is globally effective, we will compute the geometric mean of all the initial partitioning and the geometric mean of all the final results. Then, we will normalize w.r.t the first value. Doing so, we can understand whether the considered heuristic performs better or worse in a consistent way. This normalized geometric mean is denoted, in the following tables, with the symbol ρ .

4.1 Test matrices

In order to have insightful results, we mentioned that the matrices used in the numerical experiments should have different features: in particular we distinguish between rectangular matrices and square matrices, and try to have a wide selection w.r.t. the number of nonzeros.

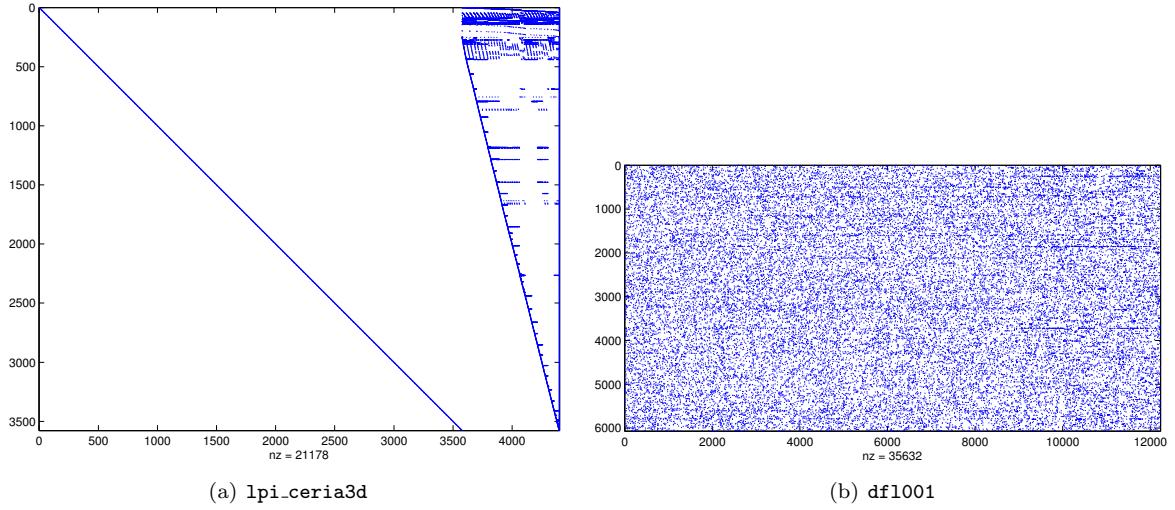
These matrices are mainly from the University of Florida Sparse Matrix Collection [24], and some can additionally be found on the Matrix Market collection [26]. The matrices `tbdmatlab` and `tbdlinux` are from [3]. Table 4.1 provides a more thorough description of the matrices used, along with an outline of their basic properties (number of rows m , number of columns n , number of nonzeros N) and their

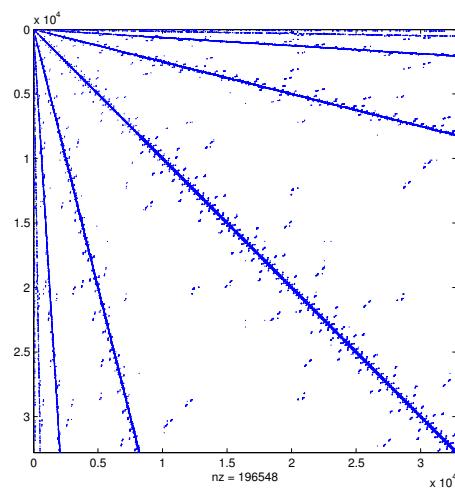
original purpose. Some of these matrices (namely the ones with the \dagger symbol in the table) belong to the 10th Dimacs Implementation Challenge [33], which addressed the graph partitioning and graph clustering problem, and are therefore naturally suited for testing the quality of the solutions produced by our algorithms. The three symmetric matrices which represent street networks, are the undirected and unweighted version of the largest strongly connected component of the corresponding OpenStreetMap road network of that country.

Name	m	n	N	Source problem
lpi_ceria3d	3576	4400	21178	Netlib Linear Programming
df1001	12230	6071	35632	Netlib Linear Programming
delaunay_n15 \dagger	32768	32768	196548	Delaunay triangulations of random points in plane
deltaX	68600	21961	247424	High fillin with exact partial pivoting
cre_b	9648	77137	260785	Netlib Linear Programming
tbdmatlab	19859	5979	430171	Term-by-document matrix
nug30	52260	379350	1567800	Netlib Linear Programming
coAuthorsCiteseer \dagger	227320	227320	1628268	Citation and coauthor network
bcsstk32 \dagger	44609	44609	2014701	Stiffness matrix for automobile chassis
bcsstk30 \dagger	28924	28924	2043492	Stiffness matrix for off-shore generator platform
wave \dagger	156317	156317	2118662	3D finite elements
tbdlinux	112757	20167	2157675	Term-by-document matrix
rgg_n_2_18_s0 \dagger	262144	262144	3094566	Random graph
belgium_osm \dagger	1441295	1441295	3099940	Street network of Belgium
polyDFT	46176	46176	3690048	Polymer self-assembly
netherlands_osm \dagger	2216688	2216688	4882476	Street network of The Netherlands
cage13	445315	445315	7479343	DNA Electrophoresis
italy_osm \dagger	6686493	6686493	14027956	Street network of Italy

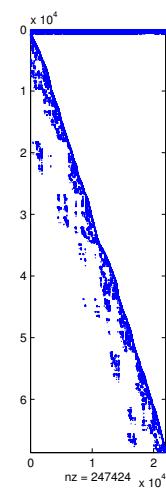
Table 4.1: Matrices used in our experiments, sorted by number of nonzeros.

In Figure 4.1 we show the sparsity patterns of the test matrices. The images were obtained using the `spy` function of MATLAB.

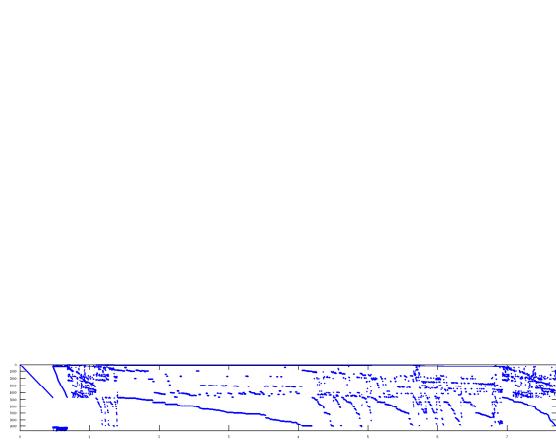




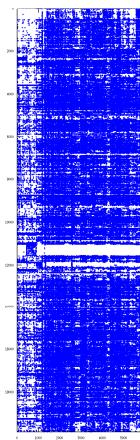
(c) `delaunay_n15`



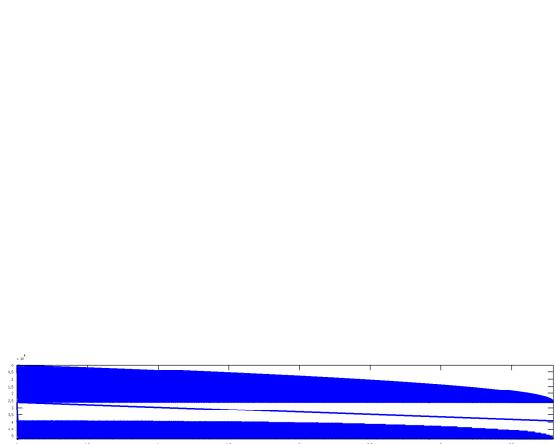
(d) `deltaX`



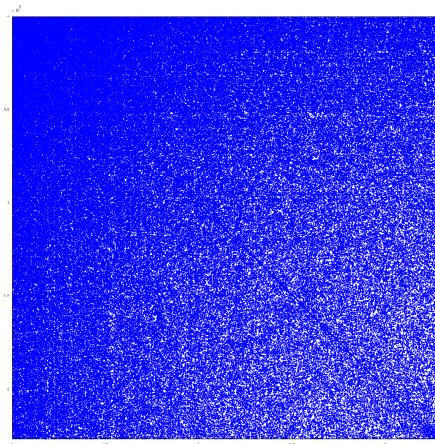
(e) `cre_b`



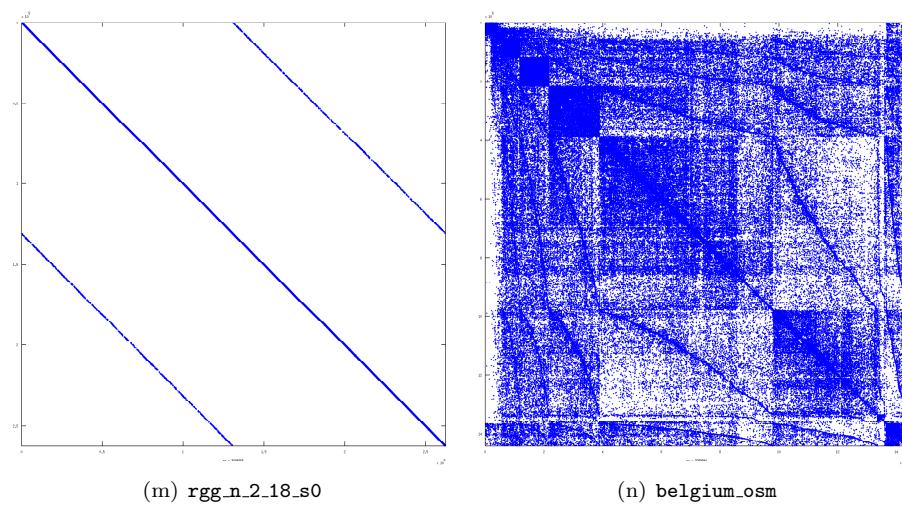
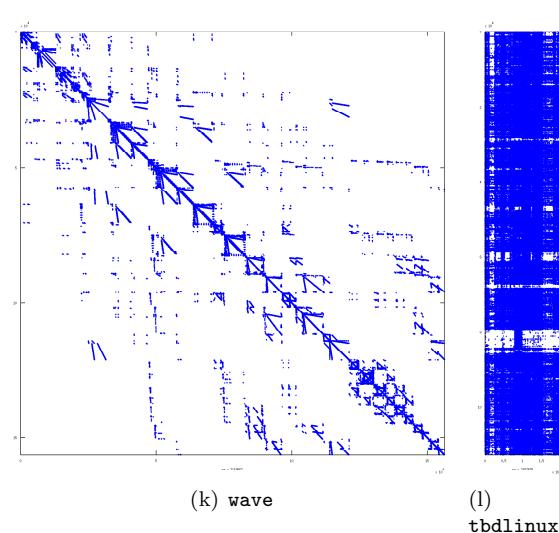
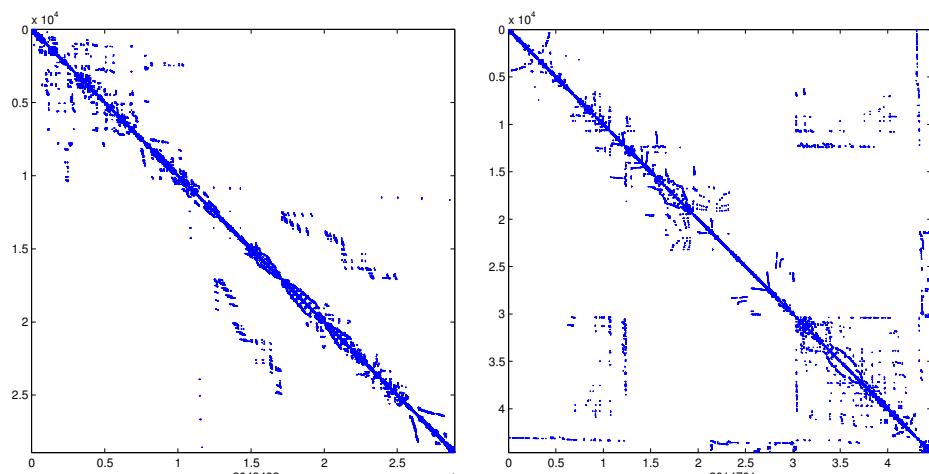
(f) `tbdmatlab`



(g) `nug30`



(h) `coAuthorsCiteSeer`



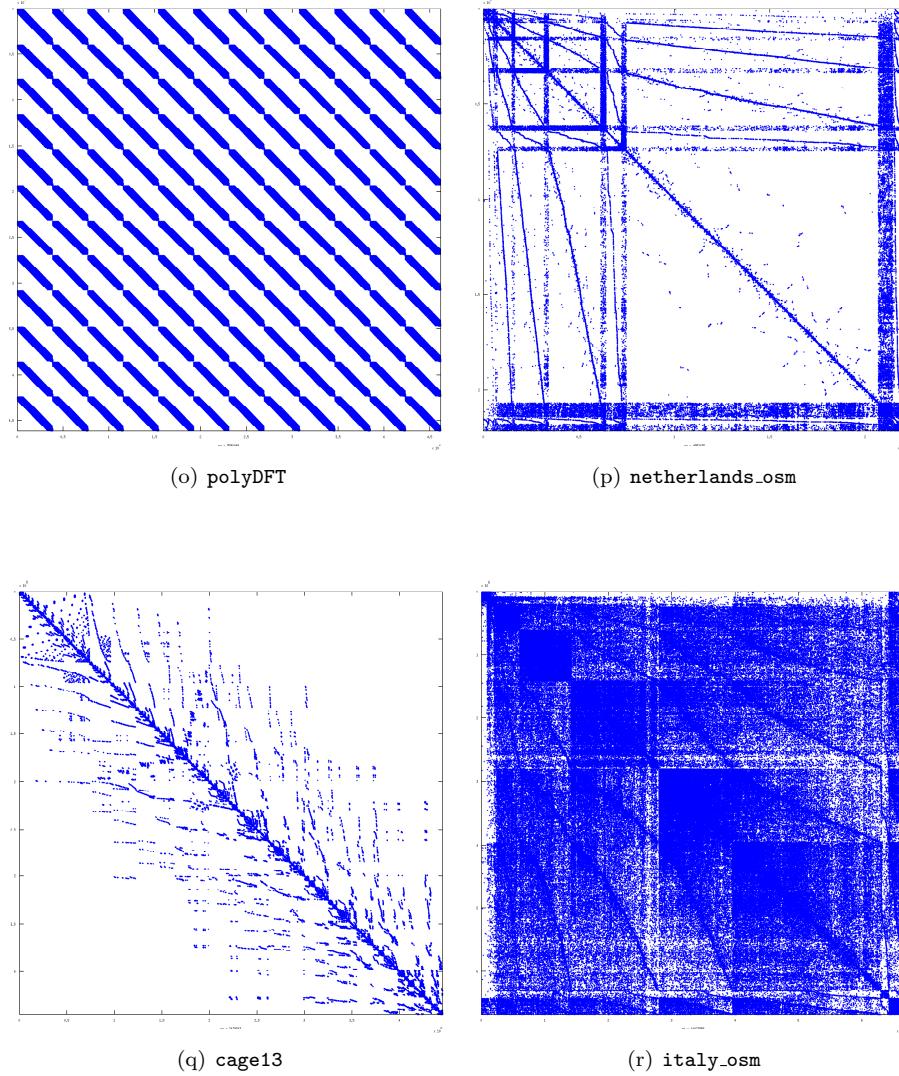


Figure 4.1: Sparsity patterns of the test matrices.

4.2 Preliminary selection of the best heuristics

In Chapter 2 and 3 we discussed many heuristics, which in turn depend on different parameters. It is best to perform a preliminary analysis to quickly figure out which heuristics produce the best solutions and should therefore be tested extensively, and the ones that should not be further considered.

For this reason, a small number of different matrices (with a relatively small number of nonzeros) has been selected from our choice of Table 4.1: in particular the considered matrices, which have a different structure, are `df1001`, `tbdlinux`, `nug30`, `rgg_n_2_18_s0`, `bcsstk30`.

The local search heuristic given in Section 2.3, quickly turned out to be far from effective, and therefore we decided to discard it altogether in the numerical experiments.

4.2.1 Partition-oblivious heuristics

Table 4.2 summarizes the results of this preliminary analysis of the partition-oblivious heuristics for the 5 chosen matrices. In the first line, the results of the medium-grain method with the algorithm proposed

in [13] are given. The value ρ represents the geometric mean of the results for a given heuristic, averaged over all matrices and normalized w.r.t. the default medium-grain method.

Heuristic	Matrix					ρ
	df1001	nug30	bcsstk30	tbdlinux	rgg_n_2_18_s0	
medium-grain	590	36262	552	8135	910	1.0
po_localview	571	36665	598	8327	1160	1.07
po_unsorted_concat_row	1492	189689	653	15081	1098	2.04
po_unsorted_concat_col	589	38491	600	24024	1066	1.32
po_unsorted_random	1314	113070	1127	20154	1093	2.11
po_unsorted_mix_alt	1461	181216	715	27942	1104	2.32
po_unsorted_mix_spr	1322	81915	759	18584	1122	1.81
po_sorted_w_simple	597	38383	785	8307	1093	1.13
po_sorted_nw_simple	606	38674	789	8301	1096	1.14
po_sorted_w_concat_row	1486	189681	642	15082	1078	2.01
po_sorted_w_concat_col	597	38655	621	24045	1068	1.33
po_sorted_nw_concat_row	1496	189683	614	15086	1090	2.01
po_sorted_nw_concat_col	593	38513	621	24005	1076	1.33
po_sorted_w_mix_alt	1317	162549	790	23683	1091	2.19
po_sorted_w_mix_spr	641	163013	782	23441	1093	1.88
po_sorted_nw_mix_alt	1457	62273	797	15015	1096	1.69
po_sorted_nw_mix_spr	719	62402	793	15072	1106	1.47
po_is	594	30655	615	13286	-	1.12

Table 4.2: Results of the devised partition-oblivious heuristics for the five chosen matrices. In each column, we use boldface to highlight the best found partitioning (not considering the medium-grain value).

In the table there is no result for the heuristic `po_is` and the matrix `rgg_n_2_18_s0`, because of memory limits.

The first method shown in the table, `po_localview`, discussed in Section 2.1, shows the best results among these partition-oblivious heuristics. This is not surprising, as such algorithm is quite similar to the one originally proposed in [13]; `po_localview` performs slightly worse (7%, on average) than `medium-grain` because the latter also includes an iterative refinement procedure [13, Section 3.3], which we will discuss in Section 4.5.

From the table, it appears that the framework discussed in Section 2.4 is not particularly effective in this case: some methods result even in twice the communication volume, on average; the `po_sorted_w_simple` and `po_sorted_nw_simple` heuristics produce the best results with a communication volume, respectively, 13% and 14% worse than the reference value. Moreover, it appears that mixing rows and columns in the priority vector is not advisable: for the matrices `df1001` and `nug30`, for example, we obtain a communication volume 3-5 times higher than `medium-grain`.

Computing the maximum independent set on the full matrix seems to be an interesting approach, especially with `nug30`: the communication volume is, on average, 16% lower than the one obtained with the medium-grain method; with the other matrices (in particular `tbdlinux`), however, the results are not as good.

Because of these preliminary results, it appears that it is interesting to test more thoroughly the heuristics `po_localview` and `po_is`, in Section 4.4.

4.2.2 Partition-aware heuristics

In Table 4.3, we summarize the results for the partition-aware heuristics. As now we are considering a fully iterative framework, it is best to explicitly give the average of the 20 initial partitionings (iteration 0) alongside the average of the 20 final results (iteration 1). The value ρ represents, similarly as before, the geometric mean of the final results, normalized w.r.t. the average of the initial partitionings.

The computation of the independent set was, once again, not possible with `rgg_n_2_18_s0`, because of memory limits.

From the table, it appears that the approaches discussed in Section 2.2 are not effective: in both cases (using the SBD and SBD2 forms of the partitioned matrix) the communication volume is, on average, more than twice the one obtained with the medium-grain method.

Regarding the framework discussed in Section 2.4, differently from the partition-oblivious case, it seems that mixing rows and columns does not produce a sharp decrease in the quality of the solutions, albeit far from being a good result: if no sorting is performed, the communication volume is on average 18% and 13% higher (depending on the mixing strategy), whereas with sorting we have a 16% and 10% worse solution.

Moreover, it appears that moving the indices with one nonzero to the back of our priority vector (denoted by `w` in the heuristics) or not (denoted by `nw`) does not yield a substantial difference; in addition, there is no clear advantage of one strategy over the other. Therefore, should we decide to consider these methods for further testing, only one of the two strategies should be picked.

The good results produced by the heuristics `pa_unsorted_concat_row` and `pa_unsorted_concat_col` (and their sorted counterparts, which are not significantly different) are a bit surprising: we expected that a more elaborate strategy (for example, sorting and mixing of rows and columns) would yield a lower communication volume than a simple concatenation. These two methods, especially for strongly rectangular matrices, produce fairly low communication volumes, if we start the concatenation with the longer dimension. If there are more rows than columns, for example, it means that the rows are in general shorter, and therefore, by giving them high priority, we have a higher chance that keeping all of their nonzeros together will not cause communication also for the columns. From another point of view, we could interpret this heuristic in the light of Chapter 3: the set of indices of the rows (or the columns) is in fact an independent set, as it is obvious that two rows (or columns) do not share any nonzero; this means that, in these two heuristics, we are roughly doing the same as in Section 3.3, albeit not as targeted to cut and uncut indices.

The behavior of these two heuristics is also consistent with the size of the matrix: if one produces good results, the other performs much worse. This suggest us that these two schemes could be merged in a single method, `pa_localbest`, which tries both and picks the best. This approach is very similar to the `localbest` method (hence the name) already employed by Mondriaan, in which the next split direction is decided in this fashion.

The computation of the maximum independent set on the uncut indices, seems to be an effective strategy for a fully iterative partitioning scheme. The average of the partitionings is indeed very close to the one obtained with the medium-grain model (rounded up, only 1% higher) and we have a definite improvement for the matrix `tbdlinux`, where the communication volume is 8% lower.

In general, we can see how with these partition-aware heuristics the results are better than with the partition-oblivious ones as the value of ρ is closer to 1. Therefore, to perform an effective selection, we need to be stricter: `pa_unsorted_concat_row` and `pa_unsorted_concat_col` (for which we will drop the label `unsorted_concat`, as `row` or `col` is sufficient to distinguish from the others), the `pa_localbest` scheme defined as above, `pa_sorted_w_simple` (which we will denote simply by `pa_simple`), `pa_is_1` and `pa_is_3` are selected for further testing in Section 4.4.

Heuristic	i	Matrix					ρ
		df1001	nug30	bcsstk30	tbdlinux	rgg_n_2_18_s0	
pa_localview	0	582	36224	537	8051	914	1.26
	1	575	36896	577	9934	2189	
pa_sbdview	0	590	36057	546	8112	899	2.17
	1	1493	187241	699	19852	1074	
pa_sbd2view	0	583	36123	542	8018	906	2.17
	1	1276	125009	1150	20757	1055	
pa_unsorted_concat_row	0	584	36512	597	7934	929	1.07
	1	628	42581	641	7337	1088	
pa_unsorted_concat_col	0	589	35945	574	7999	898	1.11
	1	589	38862	602	10087	1069	
pa_unsorted_random	0	591	36390	550	8044	909	1.08
	1	622	36952	589	8383	1094	
pa_unsorted_mix_alt	0	592	36097	536	8019	896	1.18
	1	633	40286	684	9683	1108	
pa_unsorted_mix_spr	0	589	36137	542	8024	905	1.13
	1	642	39703	645	8361	1094	
pa_sorted_w_simple	0	587	36498	566	8018	903	1.06
	1	586	38730	577	7992	1102	
pa_sorted_nw_simple	0	596	35837	568	8014	882	1.06
	1	594	38683	578	8004	1100	
pa_sorted_w_concat_row	0	585	36158	556	8019	880	1.10
	1	621	44930	617	7349	1071	
pa_sorted_w_concat_col	0	595	36157	545	7995	908	1.13
	1	596	38823	633	10108	1075	
pa_sorted_nw_concat_row	0	598	36661	548	8025	889	1.08
	1	638	42580	601	7359	1066	
pa_sorted_nw_concat_col	0	580	36421	561	8035	916	1.12
	1	595	38614	617	10084	1077	
pa_sorted_w_mix_alt	0	586	36566	549	8006	912	1.16
	1	692	41963	593	9476	1105	
pa_sorted_w_mix_spr	0	593	36085	537	8010	924	1.10
	1	619	39697	571	9060	1078	
pa_sorted_nw_mix_alt	0	588	36509	546	8020	891	1.15
	1	687	42542	566	9486	1085	
pa_sorted_nw_mix_spr	0	589	36197	553	8006	886	1.13
	1	655	39201	602	9076	1094	
pa_is_1	0	592	36716	534	7997	-	1.01
	1	588	36118	614	7323	-	
pa_is_2	0	596	35972	539	8011	-	1.06
	1	592	36375	631	8681	-	
pa_is_3	0	590	36432	540	7997	-	1.02
	1	589	36324	635	7410	-	

Table 4.3: Results for partition-oblivious heuristics. Boldface is used to highlight the best found partitioning for each matrix. Iteration 0 corresponds to the average of the initial partitionings, whereas iteration 1 is the average of the final partitionings. In the last column, for each method, the geometric mean of iteration 1 is normalized w.r.t. the geometric mean of iteration 0 for the same method.

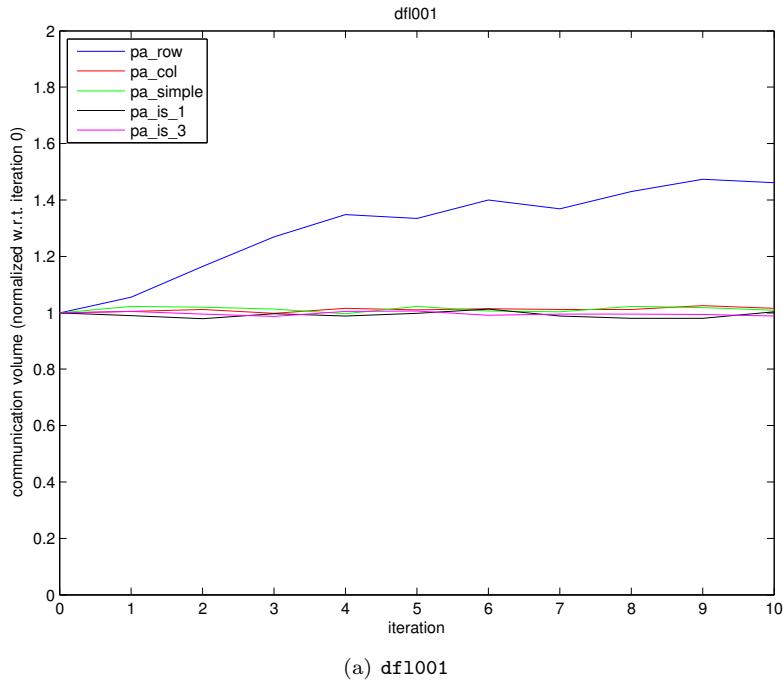
4.3 Number of iterations

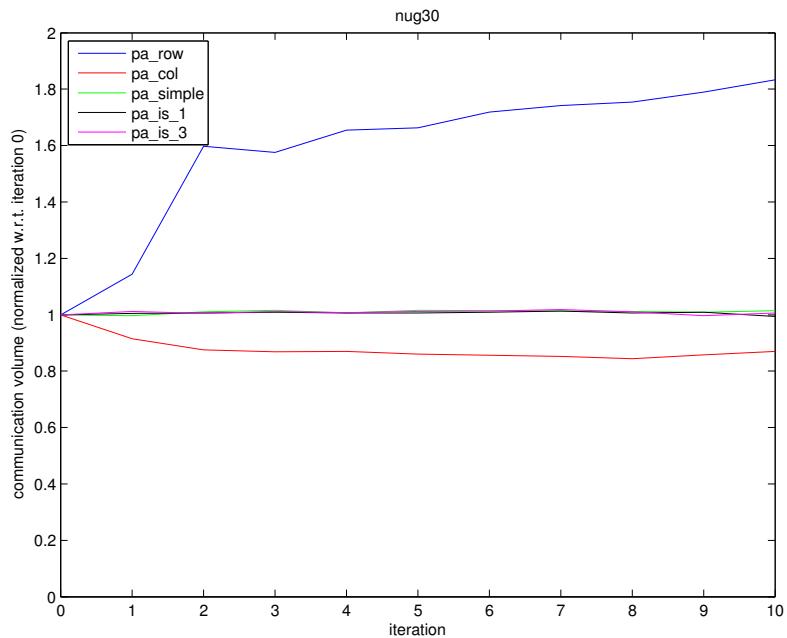
The parameter $iter_{max}$ of Algorithm 4.2 is of vital importance for the running time of our iterative method: our research motivation, for the partition-aware heuristics, is to trade computation time for solution quality, but efficiency is also a fundamental goal. An iterative scheme in which each iteration is only a slight improvement on the previous one (which means a higher value of $iter_{max}$ is needed to attain a low communication value), is less desirable than a scheme that requires just one or two iterations. Note that this reasoning only applies to partition-aware heuristics, as the partition-oblivious ones are not fully iterative: $iter_{max}$, in this case, means that a number of independent iterations is performed.

Luckily enough, the heuristics introduced in Chapters 2 and 3 are quite fast at showing their potential: whenever an improvement over the initial partitioning can be achieved, it usually happens in the first iteration; the solution quality, over the next few ones, either remains more or less constant, or gets worse. This is the main reason why in Section 4.2, for the partition-aware heuristics, we set the parameter $iter_{max}$ to 1, other than conveniently saving computation time.

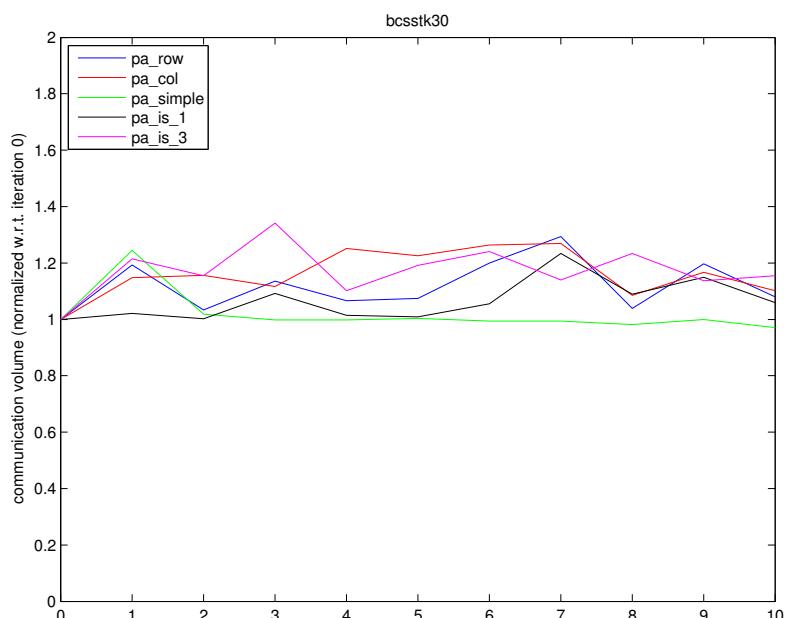
In Table 4.4, we show the results of multiple iterations for some of the selected heuristics of Section 4.2.2. We replaced the matrix `rgg_n_2_18_s0` with the matrix `delaunay_n15` because, as seen previously, the heuristics which computed the maximum independent set failed to return a solution. In Figure 4.2 these results are depicted graphically; moreover, to better convey the significance of the data, each communication volume has been normalized w.r.t. iteration 0.

As already argued in Section 4.2.2, we simplify the notation for the considered heuristics: `pa_row` and `pa_col` instead, respectively, of `pa_unsorted_concat_row` and `pa_unsorted_concat_col`, and likewise we use `pa_simple`, instead of `pa_sorted_w_simple`.

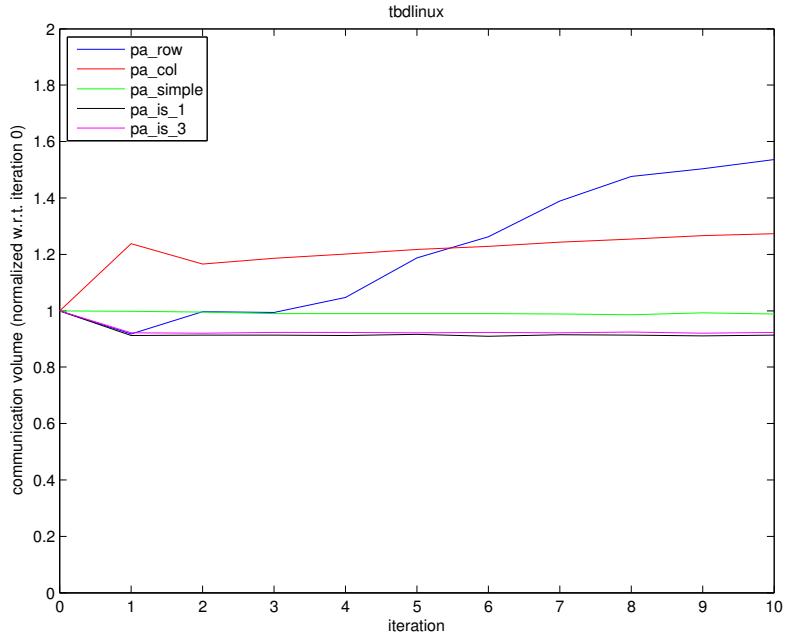




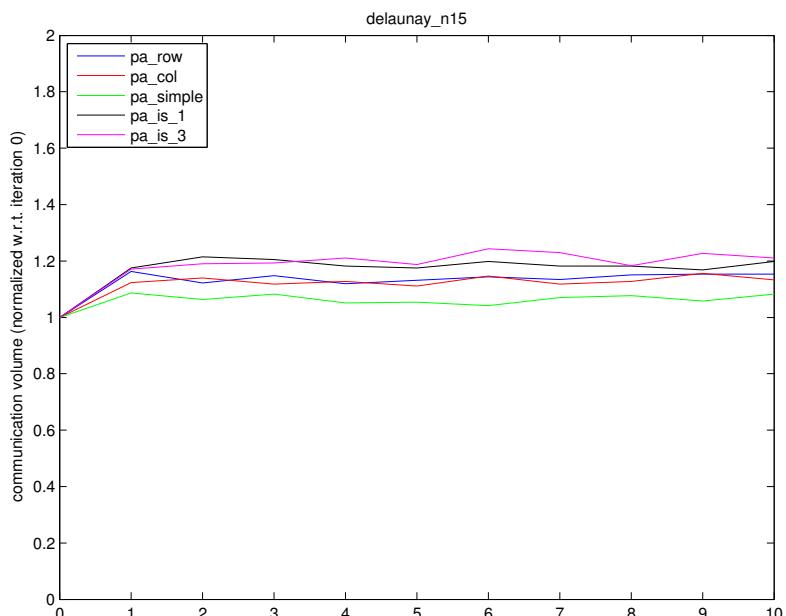
(b) nug30



(c) bcsstk30



(d) tbdlinux



(e) delaunay_n15

Figure 4.2: Graphical visualization of Table 4.4. Each value has been normalized w.r.t. the first iteration.

It seems there is no clear advantage in performing 10 iterations over just one, except for the combination `pa_unsorted_concat_col` and `nug30`: even in this case, after an improvement over the first two iterations, the communication value stagnates around 30500-31000.

In most of the other cases, it appears that there is no advantage in performing 10 consecutive iterations over just one: improvements are found immediately, or not at all.

Matrix	Heuristic	Iterations										
		0	1	2	3	4	5	6	7	8	9	10
df1001	pa_row	598	631	696	759	806	798	837	818	855	881	874
	pa_col	587	590	594	586	596	593	595	594	594	602	596
	pa_simple	589	602	601	597	587	602	593	591	602	600	594
	pa_is_1	599	593	586	597	592	598	607	592	587	587	601
	pa_is_3	597	600	594	589	600	601	592	594	594	593	590
nug30	pa_row	36308	41522	58016	57194	60054	60361	62408	63244	63699	64994	66534
	pa_col	36059	32965	31550	31335	31390	31009	30881	30709	30437	30924	31345
	pa_simple	37158	37051	37515	37723	37389	37732	37683	37660	37613	37491	37681
	pa_is_1	36649	36828	36849	36976	36856	36894	36970	37105	36851	36989	36440
	pa_is_3	35793	36191	35975	36190	36054	36219	36235	36432	36159	35671	36020
bcsstk30	pa_row	562	670	581	638	599	604	674	727	584	673	607
	pa_col	540	620	624	603	676	662	682	685	586	630	595
	pa_simple	544	677	554	543	543	546	541	541	534	544	528
	pa_is_1	562	574	563	614	570	567	593	693	612	646	595
	pa_is_3	542	658	626	727	597	646	672	618	669	616	626
tbdlinux	pa_row	8023	7365	7996	7977	8402	9527	10127	11141	11846	12059	12327
	pa_col	8108	10040	9447	9614	9732	9872	9955	10079	10171	10268	10327
	pa_simple	8044	8027	8002	7964	7958	7964	7957	7951	7934	7979	7956
	pa_is_1	8055	7343	7361	7360	7343	7375	7330	7365	7357	7331	7358
	pa_is_3	8069	7432	7426	7446	7445	7432	7444	7440	7461	7431	7443
delaunay_n15	pa_row	312	363	350	358	349	353	357	354	359	360	360
	pa_col	315	354	359	352	355	350	361	352	355	364	357
	pa_simple	313	340	333	339	329	330	326	335	337	331	339
	pa_is_1	308	362	374	371	364	362	369	364	364	360	369
	pa_is_3	305	357	363	364	369	362	379	375	361	374	369

Table 4.4: Multiple iterations of the five selected partition-aware heuristics. The numbers shown are the rounded arithmetic means of each iteration over 10 repeats of the experiment. Iteration 0 stands for the initial partitioning obtained with the medium-grain model.

4.4 Analysis of the performance of the best heuristics

In Section 4.2 we selected the most interesting heuristics for the final round of tests. In this section, we will check the quality of these methods for all of the matrices in our test bed. Our testing methodology is the same as before: 20 initial partitionings and, for each, 5 runs of a single iteration of the considered heuristic. The means are computed as in Section 4.2.

4.4.1 Partition-oblivious heuristics

The results for the partition-oblivious heuristics can be found in Table 4.5.

Matrix	Heuristic		
	medium-grain	po_localview	po_is
lpi_ceria3d	220	239	1030
df1001	590	571	594
delaunay_n15	310	338	367
deltaX	236	222	586
cre_b	580	630	632
tbdmatlab	3951	4276	4711
nug30	36262	36665	30655
coAuthorsCiteseer	9583	10459	16217
bcsstk30	552	598	615
bcsstk32	818	993	1082
wave	4624	5660	5100
tbdlinux	8135	8327	13286
rgg_n_2_18_s0	910	1160	-
belgium_osm	248	264	263
polyDFT	3633	3701	3456*
netherlands_osm	194	204	215
cage13	45233	52109	57556
italy_osm	262	280	284
ρ	1.0	1.07	1.22

Table 4.5: Results of the selected partition-oblivious heuristics with the test matrices. For each matrix, the best found average partitioning is highlighted.

First of all, we can see that the value of ρ attained by the heuristic `po_localview` is indeed very similar to the one obtained in Section 4.2, which suggests us that the preliminary matrices were indeed a good sample of the whole test set (at least for this heuristic). Regarding `po_is`, instead, we can see that now the final value of ρ is much higher: this is to be explained from the fact that in the preliminary tests one of the matrices did not yield any result, and therefore the average was computed only on four matrices, and among those there was `nug30`, and the exceptionally good result obtained with it had a lot of weight in the average. Now that there is a wide variety of matrices, such an especially good result has much less influence, and the performance of this method with the other matrices is not particularly satisfactory.

It is interesting to note that this heuristic that employs the computation of the independent set has somewhat of an erratic behavior: in a few matrices (`nug30`, `polyDFT`) the communication value is the lowest among the three methods considered (respectively 16% and 5% lower than the second best method, `medium-grain`), whereas on other matrices (`lpi_ceria3d`, `tbdlinux`, `coAuthorsCiteseer`) the final

solution is much worse.

It is interesting to highlight a particular result of the table: we marked with the symbol * the solution attained by `po_is` for the matrix `polyDFT` because, for every initial partitioning with medium-grain, each independent iteration with the heuristic produced the same communication value of 3456, the lowest recorded in general. This suggests that either this value is a local optimum hard to escape from, or either the optimal value.

4.4.2 Partition-aware heuristics

Table 4.6 shows the final results for the partition-aware heuristics, with all the matrices from Section 4.1. Also in this case we use simplified labels for the considered heuristics, as there is no possible confusion: `pa_row` and `pa_col` represent, respectively, `pa_unsorted_concat_row` and `pa_unsorted_concat_col`, whereas `pa_simple` represents `pa_sorted_w_simple`.

From the table, it appears that the properties of the matrices used for our tests were well balanced, as the final value of ρ for the heuristics `pa_row` and `pa_col` is the same (we already observed how these two heuristics show opposite behavior). Moreover, as they are, on average, only 8% worse than `medium-grain`, we confirmed our expectations: simple concatenation of rows and columns is an interesting choice. This is even more evident with `pa_localbest`: as predicted, its average communication volume is lower than the one of its parts, being only 4% worse than `medium-grain`. A similarly good result was attained also by the other method of the same family, `pa_simple`.

The main consequence of this series of experiments is that none of the examined heuristics was found to be generally better than medium-grain: even the heuristics that employed the computation of the maximum independent set resulted in, respectively, a 5% and 8% higher communication volume.

However, the algorithms `pa_localbest`, `pa_is_1` and `pa_is_3` produced very interesting results with (strongly) rectangular matrices. In particular, if we were to consider only such matrices (i.e. `lpi_ceria3d`, `df1001`, `deltaX`, `cre_b`, `tbdmatlab`, `nug30` and `tbdlinux`) we would obtain much better results: the value of ρ for these heuristics would respectively be of 0.99, 0.99 and 0.98, showing that indeed these schemes perform better (albeit slightly) than `medium-grain`. This other alternative finding shows that the independent set approach devised in Chapter 3 is indeed worthy of deeper investigation, as we will discuss in more detail in Chapter 5.

4.5 Results with iterative refinement of the final partitionings

In Section 4.2.1 we argued that the results of the medium-grain method were better than the one we achieved with `po_localview` because, among other things, the former employed an iterative refinement procedure [13], outlined as follows.

After partitioning the matrix A into two sets, A_0 and A_1 , we can use these sets to create the matrix B of the medium-grain method: for example, $A_r := A_0$ and $A_c := A_1$. Now, we do not use this matrix B as in the medium-grain method to obtain a new partitioning, but instead we aim at improving the current one. To encode the current partitioning in this matrix B (thus retaining the same communication volume and load balance), we assign its first n columns (thus, the columns of A_c) to a single processor and likewise for the other m columns (i.e. the ones of A_r^T).

Then, the hypergraph is created from this matrix B and a single run of the Kernighan-Lin method (with the improvements of Fiduccia and Mattheyses) is performed. After a step of such refinement is performed, we can repeat the whole procedure and use the obtained partitioning to compute a new matrix B . Once the communication volume does not improve anymore, we create the matrix B by swapping the roles of A_0 and A_1 and restart the whole process. This is repeated several times (continuously swapping the roles of A_0 and A_1) until no further improvement is obtained.

Note that, as the Kernighan-Lin method is monotonically non-increasing, during this whole procedure the communication volume is either lowered or remains at the same value. These few steps are therefore

Matrix	Heuristic											
	pa_row				pa_col				pa_localbest			
	initial	final	initial	final	initial	final	initial	final	initial	final	initial	final
lpi_ceria3d	225	225	223	243	223	225	224	229	221	229	222	222
df1001	584	628	589	589	589	589	587	586	593	589	590	589
delaunay_n15	310	350	306	362	310	350	312	333	309	365	309	367
deltax	236	224	236	227	236	224	234	222	235	223	237	224
cre_b	602	1014	599	637	599	637	590	627	590	649	590	627
tbdmatlab	3884	3623	3934	4096	3884	3623	3878	3903	3928	3739	3939	3677
nug30	36512	42581	35945	38862	35945	38862	36498	38731	36716	36119	36432	36324
coAuthorsCiteseer	9534	9885	9697	10050	9534	9885	9633	9955	9523	10049	9528	14737
bcsstk30	597	641	574	602	574	602	567	577	535	614	540	636
bcsstk32	804	958	837	1022	804	958	818	944	903	1037	850	1003
wave	4737	4893	4782	4893	4782	4893	4725	5427	4689	5099	4650	5199
tbdlinux	7984	7337	7999	10086	7984	7337	8018	7992	7998	7323	7997	7410
rgg_n_2_18_s0	928	1088	897	1069	897	1069	903	1102	-	-	-	-
belgium_osm	227	258	255	266	227	258	282	251	247	248	242	260
polyDFT	3643	3616	3732	3607	3732	3607	3725	3729	3648	3508	3629	3515
netherlands_osm	204	190	175	191	204	190	191	199	194	199	192	210
cage13	45845	54155	45405	52251	45405	52251	45686	49784	45252	57117	45314	57181
italy_osm	254	278	262	278	262	278	279	284	267	277	255	283
ρ	1.0	1.08	1.0	1.08	1.0	1.04	1.0	1.04	1.0	1.05	1.0	1.08

Table 4.6: Results of the selected partition-aware heuristics with the test matrices. For each matrix, we highlighted the lowest communication volume found. The column “initial” represents the average of the initial partitionings computed with medium-grain, whereas “final” represents the average of the final partitionings for that heuristic.

an effective way of refining the current partitioning.

It is interesting then to investigate the impact of this iterative refinement procedure not only on the medium-grain method (which we use to compute the initial partitionings), but also with our algorithms.

In this perspective, Table 4.7 and 4.8 correspond to Table 4.5 and 4.6.

Matrix	Heuristic		
	medium-grain	po_localview	po_is
lpi_ceria3d	222	243	136
df1001	594	577	586
delaunay_n15	310	327	318
deltaX	236	220	216
cre_b	556	589	597
tbdmatlab	3903	4199	4718
nug30	36215	37738	29110
coAuthorsCiteseer	9501	10210	10484
bcsstk30	541	755	613
bcsstk32	864	940	1039
wave	4712	5311	4727
tbdlinux	8047	8278	13253
rgg_n_2_18_s0	902	1003	-
belgium_osm	233	249	258
polyDFT	3701	3701	3456
netherlands_osm	193	196	200
cage13	45177	51753	50831
italy_osm	279	300	283
ρ	1.0	1.07	1.06

Table 4.7: Results of the partition-oblivious heuristics with iterative refinement. For each matrix, the best found average partitioning is highlighted.

We can see that with the heuristic `po_is`, the iterative refinement procedure considerably improves the final result, as the average communication volume is only 6% higher than `medium-grain` (down from 22%).

With the other partition-oblivious heuristic, however, the results hardly improved, even if theoretically the average communication volume should have been closer to the one of `medium-grain`. This discrepancy might be explained in the fact that, in reality, in the default algorithm for the medium-grain method, the assignment of nonzero to A_r or A_c works in a slightly different way: as outlined in Algorithm 2.1, in our heuristic `po_localview`, we assign such nonzero to the shorter dimension (and break ties uniformly), whereas on `medium-grain` the assignment of the nonzeros that are alone in a row or a column, and only with those, is reversed. This is the same reasoning for the parameters `w` and `nw` in Section 2.4. In our experiments, we concluded that this operation was not particularly effective, whereas in the case of this partition-oblivious heuristics it has a great beneficial effect.

Matrix	Heuristic											
	pa_row				pa_col				pa_localbest			
	initial	final	initial	final	initial	final	initial	final	initial	final	initial	final
lpi_ceria3d	214	215	220	241	214	215	215	228	223	225	229	222
df1001	580	608	590	587	590	587	585	585	597	585	594	585
delaunay_n15	306	306	314	306	306	306	306	326	310	318	305	319
deltaX	235	217	233	215	233	215	236	219	238	220	239	220
cre_b	606	671	666	595	606	595	586	619	602	590	586	606
tbdmatlab	3930	3717	3910	3965	3930	3717	3853	3860	3914	3695	3954	3640
nug30	35824	40303	36474	36695	36474	36695	36016	39062	36146	36077	36413	36222
coAuthorsCiteseer	9620	8716	9587	8685	9587	8685	9486	9802	9421	8873	9499	10829
bcsstk30	566	591	556	577	556	577	543	553	545	611	557	623
bcsstk32	827	993	897	978	897	978	850	933	783	1044	815	1011
wave	4747	4671	4669	4732	4747	4671	4693	5164	4690	4766	4732	4756
tbdlinux	7998	7303	7997	9698	7998	7303	7979	7970	8025	7321	8044	7386
rgg_n_2_18_s0	914	897	897	893	897	893	889	906	-	-	-	-
belgium_osm	238	265	294	257	294	257	278	263	245	258	261	262
polyDFT	3619	3619	3637	3567	3637	3567	3642	3670	3639	3488	3733	3530
netherlands_osm	193	192	196	187	196	187	206	194	237	204	212	194
cage13	45275	49039	45581	49274	45275	49039	45419	49196	45389	50147	45304	50361
italy_osm	269	280	228	265	228	265	295	284	272	278	278	278
ρ	1.00	1.02	1.00	1.01	1.00	0.99	1.00	1.02	1.00	1.01	1.00	1.01

Table 4.8: Results of the selected partition-aware heuristics with iterative refinement. The best found partitioning for each matrix is highlighted. Similarly as before, the column “initial” represents the average of the initial partitionings computed with `medium-grain`, whereas “final” represents the average of the final partitionings for that heuristic.

The results shown in Table 4.8 are indeed very interesting, as it seems that the iterative refinement procedure improves considerably our final results. For almost every algorithm, the average of the final partitionings was very close to the one obtained with the medium-grain method. Furthermore, we can see that `pa_localbest` has a communication volume, on average, 1% lower than `medium-grain`.

Interestingly, comparing Table 4.8 and Table 4.6, it seems that using iterative refinement has a bigger impact in the matrices for which we had previously a less than satisfactory results. As an example, consider the heuristic `pa_localbest` and the matrices `nug30` and `tbdlinux`: without iterative refinement, the communication values were, respectively, 8% worse and 8% better than `medium-grain`; with iterative refinement, instead, the communication values are respectively 1% worse and 8% better than the medium-grain method. We can clearly see that the results which were already very good, did not improve considerably, but the overall performance of the heuristic improved greatly because we mitigated the impact of the bad results.

The same observation can be made by considering only the results for rectangular matrices: the value of ρ for `pa_localview`, `pa_is_1` and `pa_is_3` is now, respectively, of 0.98, 0.96 and 0.96. It is clear that for these matrices iterative refinement produces an improvement, although not as noticeable as with all the matrices.

Now, we can look at the results from the point of view of the proposed heuristics, rather than the matrices involved in the testing. For example, we can consider the geometric mean of the communication volume obtained for each algorithm with iterative refinement, and normalize it w.r.t. the same value obtained without iterative refinement, obtaining a value ρ' . The results of this comparison are outlined in Table 4.9.

Heuristic	ρ'
<code>po_localview</code>	0.99
<code>po_is</code>	0.79
<code>pa_row</code>	0.94
<code>pa_col</code>	0.94
<code>pa_localbest</code>	0.95
<code>pa_simple</code>	0.98
<code>pa_is_1</code>	0.97
<code>pa_is_3</code>	0.95

Table 4.9: Comparison of the geometric means of the communication volume for each heuristic, with and without iterative refinement. The values are normalized w.r.t. the first column. In the Table we distinguished between partition-oblivious and partition-aware heuristics.

Also from this table, the effectiveness of the iterative refinement procedure is clear, as all the results are lower than 1. Moreover, it is interesting to observe how the various heuristics benefit differently from such iterative refinement. To this extent, the partition-oblivious heuristics are the ones that exhibits the most variability in the improvement: `po_localview` is only marginally improved (1% better), whereas `po_is` has a communication volume 21% lower than without iterative refinement of the final partitionings.

Chapter 5

Conclusions and further developments

The goal of this thesis was firstly to investigate the opportunity of performing sparse matrix partitioning in an iterative fashion, devising a systematic approach for the selection of the information on the previous partitioning to be kept for an improvement on the quality of the solution. Secondly, our intent was also to explore the margins of improvement of the initial solution for the medium-grain model.

We were able to deal with these two research directions at the same time, applying basic principles to come up with simple schemes that could serve either purpose (with minor modifications). In the following sections, we summarize our findings and try to envision additional developments in these directions.

5.1 Improvement of the initial solution

As already outlined in Section 1.4, for the medium-grain model the initial split of the matrix A into A_r and A_c is of vital importance: our goal was to investigate whether the algorithm proposed by the authors in [13] was indeed a good choice and, if so, to what extent. To this purpose, we devised a number of different algorithms in Section 2.1, 2.4 and 3.3, which exploited different features of the underlying matrix.

Among these proposed partition-oblivious heuristics, `po_localview` (outlined in Section 2.1) and `po_is` (described in Section 3.3) showed the best results during the preliminary testing. In the final tests, the behavior of these heuristics was slightly different: without the employment of the iterative refinement procedure, briefly described in Section 4.5, `po_localview` produced the best results (7% worse than `medium-grain`), whereas with iterative refinement, `po_is` is the heuristic that performs best (6% higher communication volume than `medium-grain`).

The behavior of `po_localview` was not surprising, as we already argued how it is a less refined version of the default algorithm for the medium-grain method, whereas the good result of `po_is` confirm our expectations that the independent set formulation of Chapter 3 is indeed worthwhile. Furthermore, it is interesting to observe that `medium-grain` was the best algorithm, outperforming the partition-oblivious schemes devised in this thesis; this suggest us that it is indeed a valid algorithm to produce the initial A_r and A_c for the medium-grain method.

5.2 Iterative partitioning

It is reasonable to assume that the initial split of the matrix A into A_r and A_c can be performed more efficiently with additional information, especially if the matrix A has already been (bi)partitioned. The

main feature of a partitioning we decided to preserve was the fact that a row/column was uncut: this means that the partitioner decided, at some point in the previous iteration, that it was convenient for the nonzeros of such row/column to be assigned to the same processor. It is then reasonable to devise heuristics which tend to give a preference for these uncut rows/columns, trying to keep them uncut also in the next iteration. Naturally, as it is not always possible to do (for example because we are not directly assigning to processors, but rather to A_r and A_c), a more sophisticated approach is usually needed.

To this extent, a wide variety of heuristics has been devised: the partition-aware extension of the original algorithm, along with the heuristics that employ the Separated Block Diagonal structure of order 1 and 2 of the partitioned matrix A , quickly turned out to be far from effective. The other considered algorithms, instead, yielded more interesting results, especially with rectangular matrices.

Chapter 3 was focused on the concept of independent set, and the experimental results suggest that this is indeed an interesting approach: our best heuristics (`pa_localbest`, `pa_is_1`, `pa_is_3`) rely indeed on this concept, implicitly or explicitly. Even the simple concatenation of rows and columns in the priority vector v can be intended as a special case of the ideas of Section 3.3: the set of rows (or columns) is as a matter of fact an independent set, although probably not of maximum cardinality. The improvement of the results of this algorithm with increasing rectangularity of the matrices is also to be seen in this perspective: the more one dimension is dominant, the more that set of indices has a cardinality closer to the maximum independent set, when computed on the graph constructed according to Section 3.1.

Without iterative refinement, described in Section 4.5, none of the proposed scheme outperforms `medium-grain` for all test matrices: `pa_localbest` and `pa_simple` have a communication volume, on average, 4% higher than `medium-grain`. Considering only the rectangular matrices in our test bed, instead, the results of these heuristics were better (albeit marginally) than `medium-grain`: `pa_is_3` had an average communication volume 2% better than the default algorithm of the medium-grain method.

With the iterative refinement procedure, instead, the results are more encouraging, as the heuristic `pa_localbest` was able to outperform `medium-grain` even if slightly: the communication volume was, on average, 1% lower than the reference value. Also the other heuristics benefited greatly from this procedure, and the quality of the final partitionings is indeed very close to the one achieved by `medium-grain`. Also in the case of rectangular matrices the results are slightly better: `pa_is_1` and `pa_is_3` have a communication volume on average 4% better than the reference value.

In general, the experimental results confirmed our theoretical expectation that the independent set approach (either implicit or explicit) is indeed worthwhile.

5.3 Further research

Even though in Section 5.1 we noted how the algorithm proposed in [13] is still the most efficient for the initial split into A_r and A_c , it might still be convenient to investigate additional strategies and compare their efficiency, in order to find a better method or gain additional confidence on this algorithm.

The partition-aware heuristics proposed in this thesis, while marginally effective, are meant as a first attempt at a fully iterative approach at sparse matrix partitioning, and further research can easily be performed. It might be worthwhile, for example, to investigate more thoroughly the properties of the maximum independent set, in order to fully exploit the bidimensionality of the medium-grain model.

Even though the implementation of the Hopcroft-Karp algorithm described in Section 3.2.2 was successful in most cases (and the computation time required was, in general, reasonable), it might be interesting to produce a C implementation, to allow an eventual integration with the Mondriaan software package.

In addition, if further research on our findings with respect to rectangular matrices is successful (i.e. our results are consistently reproducible), we could add our strategy as an optional feature to the partitioner: the program, before applying any other technique, might ask the user whether he/she intends to sacrifice computation time for a better partitioning. If so, the software would try to recognize whether the considered matrix is strongly rectangular and eventually execute our approach.

Bibliography

- [1] Jan Rabaey. *Digital integrated circuits : a design perspective*. Prentice Hall, 1996. ISBN: 0131786091 (cited on page 2).
- [2] Rob H. Bisseling. *Parallel Scientific Computation: A structured approach using BSP and MPI*. Oxford University Press, 2004 (cited on pages 2, 3, 7).
- [3] Brendan Vastenhouw and Rob H. Bisseling. “A Two-Dimensional Data Distribution Method for Parallel Sparse Matrix-Vector Multiplication”. In: *SIAM Review* 47.1 (2005), pp. 67–95 (cited on pages 2, 4, 8, 10, 31, 32).
- [4] Leslie G. Valiant. “A bridging model for parallel computation”. In: *Commun. ACM* 33.8 (Aug. 1990), pp. 103–111. ISSN: 0001-0782. DOI: 10.1145/79173.79181. URL: <http://doi.acm.org/10.1145/79173.79181> (cited on page 2).
- [5] Ümit V. Çatalyürek and Cevdet Aykanat. “Hypergraph-partitioning-based decomposition for parallel sparse-matrix vector multiplication”. In: *Parallel and Distributed Systems, IEEE Transactions on* 10.7 (1999), pp. 673–693. ISSN: 1045-9219. DOI: 10.1109/71.780863 (cited on pages 5, 6, 8).
- [6] Bruce Hendrickson. “Graph partitioning and parallel solvers: Has the emperor no clothes?” In: *Solving Irregularly Structured Problems in Parallel*. Ed. by Alfonso Ferreira, José Rolim, Horst Simon, and Shang-Hua Teng. Vol. 1457. Lecture Notes in Computer Science. Springer Berlin Heidelberg, 1998, pp. 218–225. ISBN: 978-3-540-64809-3. DOI: 10.1007/BFb0018541. URL: <http://dx.doi.org/10.1007/BFb0018541> (cited on page 6).
- [7] Sivasankaran Rajamanickam and Erik G. Boman. “Parallel partitioning with Zoltan: Is hypergraph partitioning worth it?” In: *Graph Partitioning and Graph Clustering*. AMS, 2012, pp. 37–52 (cited on page 6).
- [8] Bruce Hendrickson and Tamara G. Kolda. “Graph partitioning models for parallel computing”. In: *Parallel computing* 26.12 (2000), pp. 1519–1534 (cited on page 6).
- [9] George Karypis, Rajat Aggarwal, Vipin Kumar, and Shashi Shekhar. “Multilevel hypergraph partitioning: applications in VLSI domain”. In: *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* 7.1 (1999), pp. 69–79. ISSN: 1063-8210. DOI: 10.1109/92.748202 (cited on page 6).
- [10] David A. Papa and Igor L. Markov. “Hypergraph Partitioning and Clustering”. In: *Approximation Algorithms and Metaheuristics*. 2007 (cited on page 6).
- [11] Thomas Lengauer. *Combinatorial algorithms for integrated circuit layout*. New York, NY, USA: John Wiley & Sons, Inc., 1990. ISBN: 0-471-92838-0 (cited on page 7).
- [12] Ümit V. Çatalyürek and Cevdet Aykanat. “A fine-grain hypergraph model for 2D decomposition of sparse matrices”. In: *Proceedings of the 15th International Parallel and Distributed Processing Symposium*. 2001, pp. 609–625 (cited on pages 7, 8).
- [13] Daniël M. Pelt and Rob H. Bisseling. “A Medium-Grain Method for Fast 2D Bipartitioning of Sparse Matrices”. 2013 (cited on pages 8–11, 37, 45, 50, 51).
- [14] Bora Uçar and Cevdet Aykanat. “Revisiting hypergraph models for sparse matrix partitioning”. In: *SIAM review* 49.4 (2007), pp. 595–603 (cited on page 8).
- [15] Ümit V. Çatalyürek and Cevdet Aykanat. “A hypergraph-partitioning approach for coarse-grain decomposition”. In: *Supercomputing, ACM/IEEE 2001 Conference*. IEEE. 2001, pp. 42–42 (cited on page 8).

- [16] Yifan F. Hu, Kevin C.F. Maguire, and Richard J. Blake. “A multilevel unsymmetric matrix ordering algorithm for parallel process simulation”. In: *Computers & Chemical Engineering* 23.11 (2000), pp. 1631–1647 (cited on page 8).
- [17] Ümit V. Çatalyürek, Cevdet Aykanat, and Bora Uçar. “On two-dimensional sparse matrix partitioning: Models, methods, and a recipe”. In: *SIAM Journal on Scientific Computing* 32.2 (2010), pp. 656–683 (cited on page 8).
- [18] Ümit V. Çatalyürek and Cevdet Aykanat. “PaToH (Partitioning Tool for Hypergraphs)”. In: *Encyclopedia of Parallel Computing*. Springer, 2011, pp. 1479–1487 (cited on page 8).
- [19] George Karypis and Vipin Kumar. “Multilevel k-way hypergraph partitioning”. In: *VLSI design* 11.3 (2000), pp. 285–300 (cited on page 8).
- [20] Karen D. Devine, Erik G. Boman, Robert T. Heaphy, Rob H. Bisseling, and Ümit V. Çatalyürek. “Parallel hypergraph partitioning for scientific computing”. In: *Proceedings of the 20th international conference on Parallel and distributed processing*. IPDPS’06. Rhodes Island, Greece: IEEE Computer Society, 2006, pp. 124–124. ISBN: 1-4244-0054-6. URL: <http://dl.acm.org/citation.cfm?id=1898953.1899056> (cited on page 8).
- [21] Brian W. Kernighan and Shen Lin. “An efficient heuristic procedure for partitioning graphs”. In: *Bell system technical journal* 49 (1970), pp. 291–307 (cited on page 8).
- [22] Charles M. Fiduccia and Robert M. Mattheyses. “A linear-time heuristic for improving network partitions”. In: *Design Automation, 1982. 19th Conference on*. IEEE. 1982, pp. 175–181 (cited on page 9).
- [23] Albert-Jan N. Yzelman and Rob H. Bisseling. “Cache-oblivious sparse matrix-vector multiplication by using sparse matrix partitioning methods”. In: *SIAM Journal on Scientific Computing* 31.4 (2009), pp. 3128–3154 (cited on pages 13, 15).
- [24] Tim A. Davis and Yifan Hu. “The University of Florida Sparse Matrix Collection”. In: *ACM Transactions on Mathematical Software* 38 (2011), pp. 1–25. URL: <http://www.cise.ufl.edu/research/sparse/matrices> (cited on pages 15, 17, 32).
- [25] B. O. Fagginger Auer. “GPU Acceleration of Graph Matching, Clustering, and Partitioning”. Ph. D. thesis. Utrecht University, 2013 (cited on page 15).
- [26] Ronald F. Boisvert, Roldan Pozo, Karin A. Remington, Richard F. Barrett, and Jack Dongarra. “Matrix Market: a web resource for test matrix collections.” In: *Quality of Numerical Software*. 1996, pp. 125–137. URL: <http://math.nist.gov/MatrixMarket/> (cited on pages 17, 32).
- [27] Bruce Hendrickson and Tamara G. Kolda. “Partitioning rectangular and structurally unsymmetric sparse matrices for parallel processing”. In: *SIAM Journal on Scientific Computing* 21.6 (2000), pp. 2048–2072 (cited on page 24).
- [28] Michael R. Garey and David S. Johnson. *Computers and Intractability; A Guide to the Theory of NP-Completeness*. New York, NY, USA: W. H. Freeman & Co., 1990. ISBN: 0716710455 (cited on pages 25, 26).
- [29] Dénes König. “Gráfok és mátrixok”. In: *Matematikai és Fizikai Lapok* 38 (1931), pp. 116–119 (cited on page 26).
- [30] John E. Hopcroft and Richard M. Karp. “An $n^{5/2}$ algorithm for maximum matchings in bipartite graphs”. In: *SIAM Journal on computing* 2.4 (1973), pp. 225–231 (cited on page 28).
- [31] MathWorks. *C/C++ Source MEX-Files*. 2013. URL: http://www.mathworks.nl/help/matlab/matlab_external/c-c-source-mex-files.html (cited on page 32).
- [32] David Eppstein. *Hopcroft-Karp bipartite matching (Python recipe)*. 2002. URL: <http://code.activestate.com/recipes/123641-hopcroft-karp-bipartite-matching/> (cited on page 32).
- [33] David A. Bader, Henning Meyerhenke, Peter Sanders, and Dorothea Wagner. *10th Dimacs Implementation Challenge*. 2012. URL: <http://www.cc.gatech.edu/dimacs10/index.shtml> (cited on page 33).