

Contents

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 model	9
2	Strategies for splitting a matrix A into A_r and A_c	12
2.1	Individual assignment of nonzeros	12
2.2	Assignment of blocks of nonzeros	13
2.2.1	Using the Separated Block Diagonal form of A	14
2.2.2	Using the Separated Block Diagonal form of order 2 of A	16
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	21
3	Independent set formulation of the partial row/column assignment problem	25
3.1	Graph construction	25
3.2	Computation of the independent set	26
3.2.1	Hopcroft-Karp algorithm	26
4	Experimental results	27
4.1	Improving the initial partitioning	27
4.2	Fully iterative partitioning	27
	Bibliography	28

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 very large size and involve **sparse** matrices, because the variables at hand are usually connected to a limited number of others (for example, a very 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 $\mathcal{O}(n^2) = \mathcal{O}(mn)$ operations; this is not very efficient if we have 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 we have that $N \ll mn$.

As mentioned, the systems considered are very 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 a very 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 very 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 A_0, \dots, A_{p-1} disjoint subsets. Moreover, also the input vector v and the output vector u are distributed among the p processors (note that their distribution might not necessarily, and usually it is not, the same).

Figure 1.1 shows a possible partitioning of a 9×9 matrix with 18 nonzeros. As the 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.

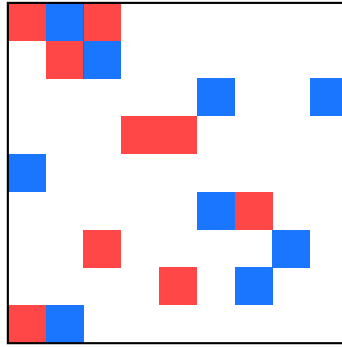


Figure 1.1: Example of a distribution among two processors of a 9×9 matrix with 18 nonzeros. Only the sparsity pattern is shown.

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 very small sum with negligible computational cost and for this reason it has been omitted from the algorithm.

Note that we assume that all of the nonzero values are all 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 the cost of a parallel sparse matrix-vector multiplication.

Figure 1.2 shows an example of the communication involved in supersteps (0) and (2), with the example partitioning shoed 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 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

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

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

```

(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.

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

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 up 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 partitioning 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.

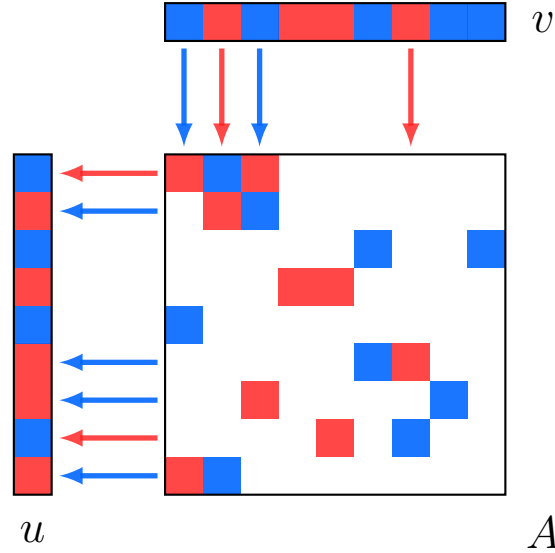


Figure 1.2: Communication for the sparse parallel matrix-vector multiplication with a matrix partitioned as in Figure 1.1. Vertical arrows represent step (0) while horizontal ones represent step (2). The color of an arrow denotes which processor is sending their data for that row/column.

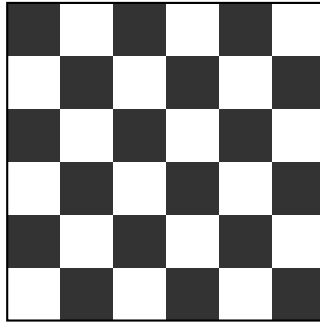


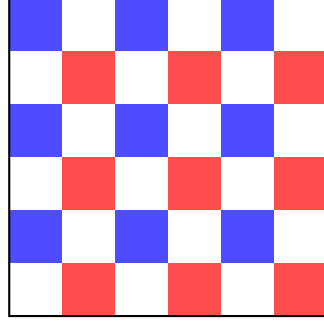
Figure 1.3: Example matrix with checkered sparsity pattern. Black boxes represent the nonzeros.

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, we have that communication is required if and only if the rows/columns of the matrices are *cut*, i.e. assigned to more than one processor.

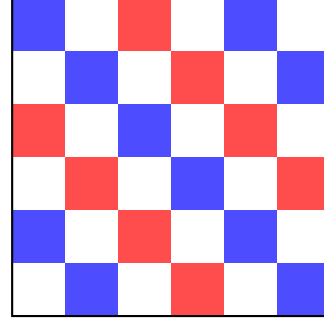
If a full column of our matrix A is assigned to the same processor, we can freely assign the corresponding component of v to the same processor, eliminating completely one source of communication (namely, the fan-out for that column). The same reasoning can be done for the rows. This simplification is possible because imposing a hard constraint similar to (1.3) also to the vector distribution is not very 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 communications cost, 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, we have that $\lambda_i, \mu_j \geq 0$.

Then the total time costs for the communication steps in our Algorithm 1.1 are:



(a) Rows and columns are not split, therefore there is no need for communication.



(b) Every row and column is split and causes communication during fan-in and fan-out.

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

$$T_{(0)} = \sum_{j=0}^{n-1} (\mu_j - 1),$$

$$T_{(2)} = \sum_{i=0}^{m-1} (\lambda_i - 1).$$
(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. Adding these costs together, we define the **communication volume** V of the considered partitioning as

$$V := V(A_0, \dots, A_{p-1}) = T_{(0)} + T_{(2)} = \sum_{i=0}^{m-1} (\lambda_i - 1) + \sum_{j=0}^{n-1} (\mu_j - 1).$$
(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 partitioning 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 *cutsizes*, 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][6], that this cut-edge metric is not an accurate representation of the communication volume. Additional criticism [7] comes from the fact that the graph partitioning approach can only handle square symmetric matrices. 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 designing 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 such a hypergraph. Each colored set represents a different hyperedge; we can see that we can have hyperedges which contain only one node.

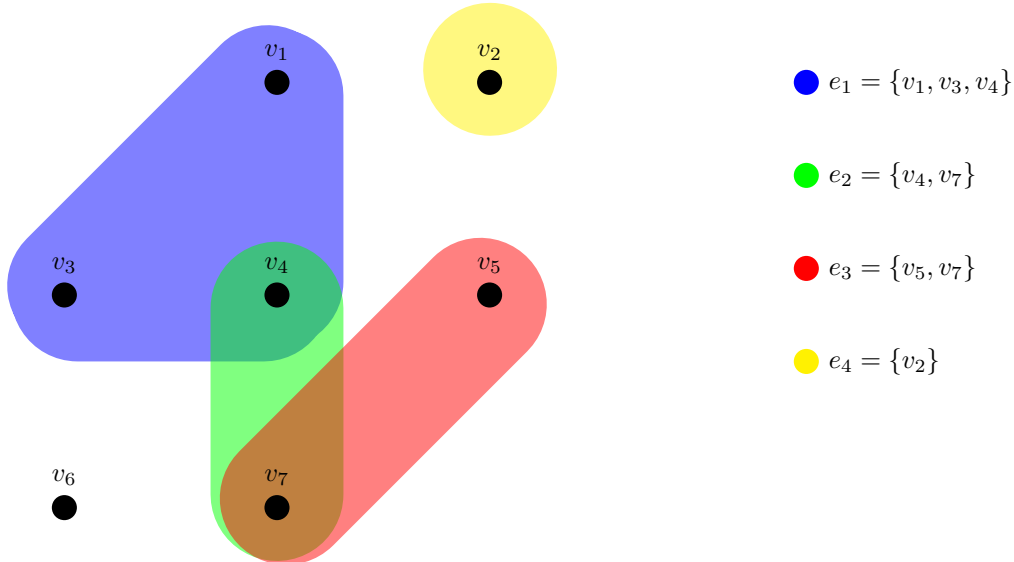


Figure 1.5: Example of a hypergraph with 7 nodes and 4 hyperedges.

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, we have that 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 partitioning hypergraph model, along with the simple 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, we have that the size of the hyperedge e_j is exactly the number of nonzeros in that column, whereas the node v_i belongs exactly to as many hyperedges as there are nonzeros in that row.

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 very 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, we have that 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. New models are relatively rare, and recently Pelt and Bisseling proposed the interesting **medium-grain** model [13]. As this model is at the very base of our work, a more detailed explanation will be given in Section 1.4.

In addition to these models, there has been some research effort towards the creation of more complicated methods, which often comprise several stages and combine different models.

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, Vastenhouw 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: the hypergraph-based approach was indeed initially devised for structurally symmetric matrices [5]. Moreover, Hu, Maguire and Blake present in [16] an algorithm for nonsymmetric matrices that performs row and column permutations, getting 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 model, 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 very different results, with respect to both solution quality and execution time, despite having at the core the same method for finding good initial solutions. The

method employed 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 very 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 with direct splitting into the desired number of parts.

1.4 Medium-grain model

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 very often too big to manage.

A new model 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 model 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 we have 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

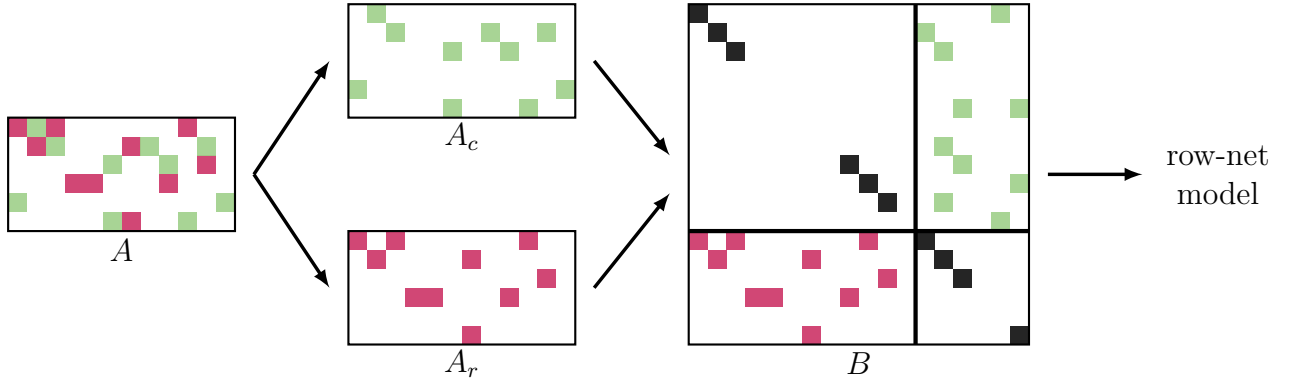


Figure 1.6: Example of the construction of the matrix B from a 6×12 matrix A , for which the sets A_r and A_c were previously established and colored differently. In the resulting matrix, the dummy nonzeros are depicted in black.

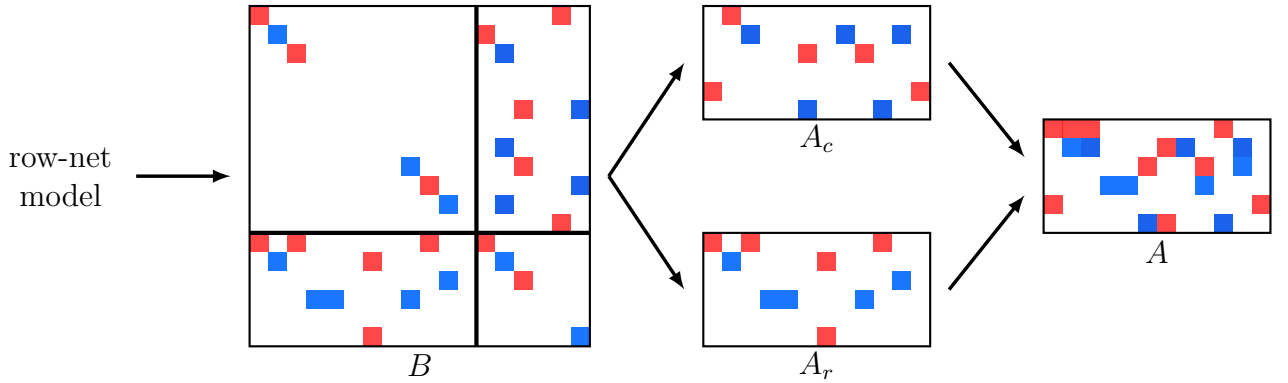


Figure 1.7: Process of obtaining a matrix partitioning starting from a partitioning of the hypergraph following the medium-grain model. In this case $p = 2$.

column of A_c , that belongs to 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 seems to confirm that this model 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 model, following two possible directions.

First of all, as the outcome of the medium-grain model 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 model, where a full multi-level partitioning is performed at each iteration and computation time is traded for solution quality.

Both these research directions share a very 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 model.

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 don't 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, despite it is impossible to do it all the time (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 have 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 and 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's exactly the same as Algorithm 1 originally proposed.

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 between  $c$  and  $r$ 
end if
 $A_r := A_c := \emptyset$ 
for all  $a_{ij} \in A$  do
    if  $nz(i) < nz(j)$  then
        assign  $a_{ij}$  to  $A_r$ 
    else if  $nz(j) < nz(i)$  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:

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 trying to assign at the same time a greater amount of nonzeros to either A_r or A_c . In particular, we will discuss how to exploit the Separated Block Diagonals (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 only 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 we have, for example, that $A([1, \dots, m], [1 \dots n]) = A$. 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 between  $c$  and  $r$ 
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  $nz(i) < nz(j)$  then
            assign  $a_{ij}$  to  $A_r$ 
        else if  $nz(j) < nz(i)$  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 $i = 0, 1, 2$, we have that the SBD form is the resulting block matrix:

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

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;

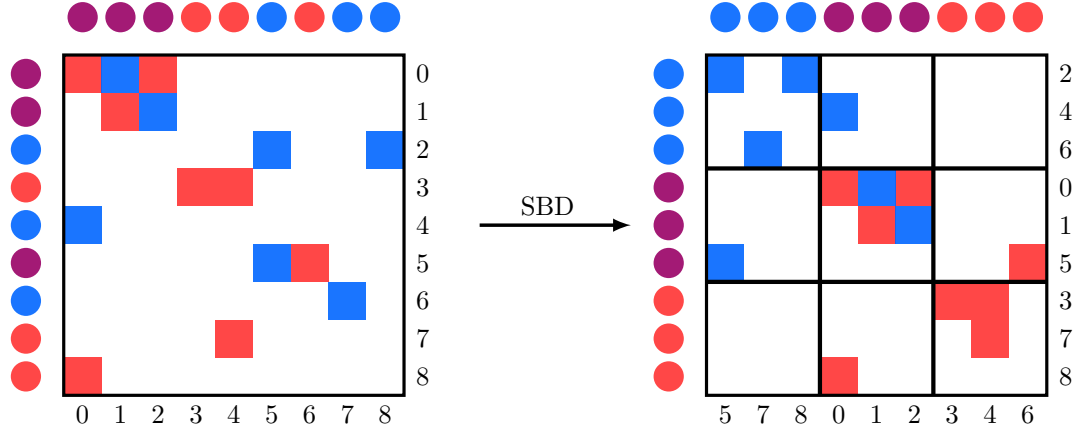


Figure 2.1: Example process to obtain the SBD form of a partitioned matrix. On the left the original matrix is shown, whereas on the right we have the permuted SBD form. On the top/left sides of the matrices the color of the circle denotes whether that row/column is completely red or blue or it is mixed (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 amount of contained nonzero 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 of \dot{A}_{00} and \dot{A}_{22} ; conversely, if the matrix has a very dense (or complicated) pattern and/or the partitioning is far from the 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: a convenient thing is to base our choice on the sizes of such blocks. For example, if $m_0 < n_0$, in the block \dot{A}_{00} we have that the columns are (on average) longer than the rows: if we assign the nonzeros of this block to A_c we are, on principle, making sure that more things 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 can’t 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 is described explicitly in Algorithm 2.3.

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

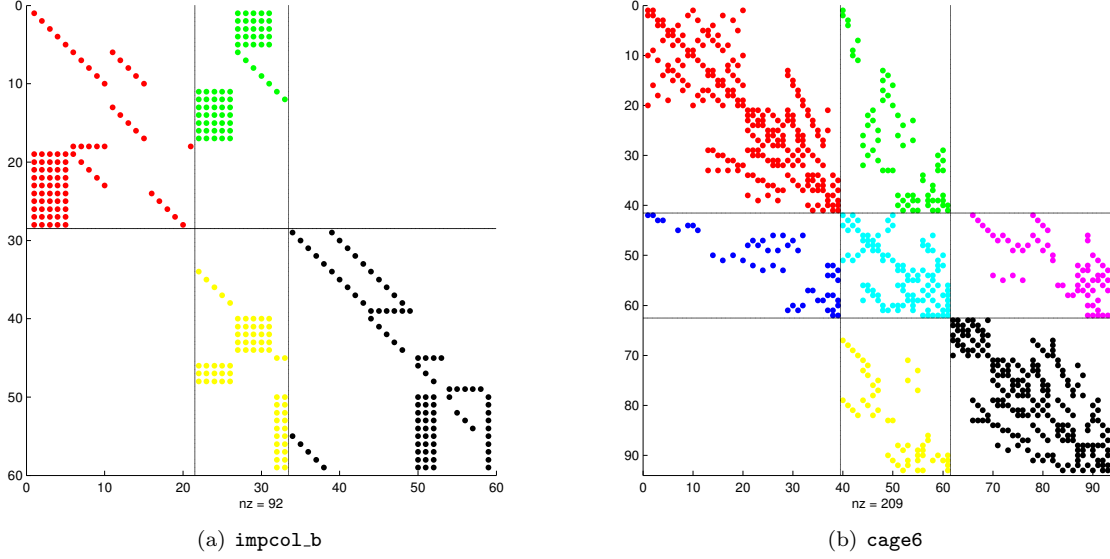


Figure 2.2: Example of SBD forms of partitioning of the matrices `impcol_b` [24] and `cage6` [25]. Each part of \dot{A} has been colored differently. In the first matrix there are no cut rows, which means that $\dot{A}_{10} = \dot{A}_{11} = \dot{A}_{12} = \emptyset$.

be easily used to compute the SBD form of a matrix [23], instead of computing this form from scratch as we described earlier.

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}_{20} the empty and non-empty columns, and in \dot{A}_{01} and \dot{A}_{02} the empty and non-empty rows. Then all the other blocks, except the central one, are permuted and split up accordingly. This procedure is better shown in Algorithm 2.4.

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

$$\ddot{A} := \begin{bmatrix} \ddot{A}_{00} & \ddot{A}_{01} & & & \\ \ddot{A}_{10} & \ddot{A}_{11} & \ddot{A}_{12} & & \\ & \ddot{A}_{21} & \ddot{A}_{22} & \ddot{A}_{23} & \\ & & \ddot{A}_{32} & \ddot{A}_{33} & \ddot{A}_{34} \\ & & & \ddot{A}_{43} & \ddot{A}_{44} \end{bmatrix}, \quad (2.2)$$

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.

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 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 \dot{A}_{01} don't have any neighbor (w.r.t. their row) with cut columns but have neighbors (w.r.t.

Input: partitioned matrix A

Output: A_r, A_c

```

 $\dot{A} :=$  SBD form of the partitioned matrix  $A$ .
 $m_0 := |\{i : \text{row } i \text{ is fully assigned to processor 0}\}|$ 
 $m_2 := |\{i : \text{row } i \text{ is fully assigned to processor 1}\}|$ 
 $n_0 := |\{j : \text{column } j \text{ is fully assigned to processor 0}\}|$ 
 $n_2 := |\{j : \text{column } j \text{ is fully assigned to processor 1}\}|$ 
 $A_r := A_c := \emptyset$ 
if  $m_0 < n_0$  then
    Assign nonzeros of  $\dot{A}_{00}$  to  $A_r$ 
else
    Assign nonzeros of  $\dot{A}_{00}$  to  $A_c$ 
end if
if  $m_2 < n_2$  then
    Assign nonzeros of  $\dot{A}_{22}$  to  $A_r$ 
else
    Assign nonzeros of  $\dot{A}_{22}$  to  $A_c$ 
end if
Assign nonzeros of  $\dot{A}_{11}$  to  $A_r$  or  $A_c$  following Algorithm 2.1
Assign nonzeros of  $\dot{A}_{10}$  to  $A_c$ 
Assign nonzeros of  $\dot{A}_{12}$  to  $A_c$ 
Assign nonzeros of  $\dot{A}_{01}$  to  $A_r$ 
Assign nonzeros of  $\dot{A}_{21}$  to  $A_r$ 

```

Algorithm 2.3: Assignment of the nonzeros based on the SBD form of the partitioned matrix 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 2.4: Algorithm to obtain SBD2 form of a matrix A .

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.

These properties are interesting because, now, nonzeros in \ddot{A}_{00} and \ddot{A}_{44} can be entirely removed from the original partitioning problem: they constitute a subset of the nonzeros that can be *perfectly* partitioned, without causing any communication. The size of these parts, and more generally of all of the blocks of \ddot{A} , is again highly dependant on the structure of the matrix, as it is shown in Figure 2.4.

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 in Algorithm 2.4.

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

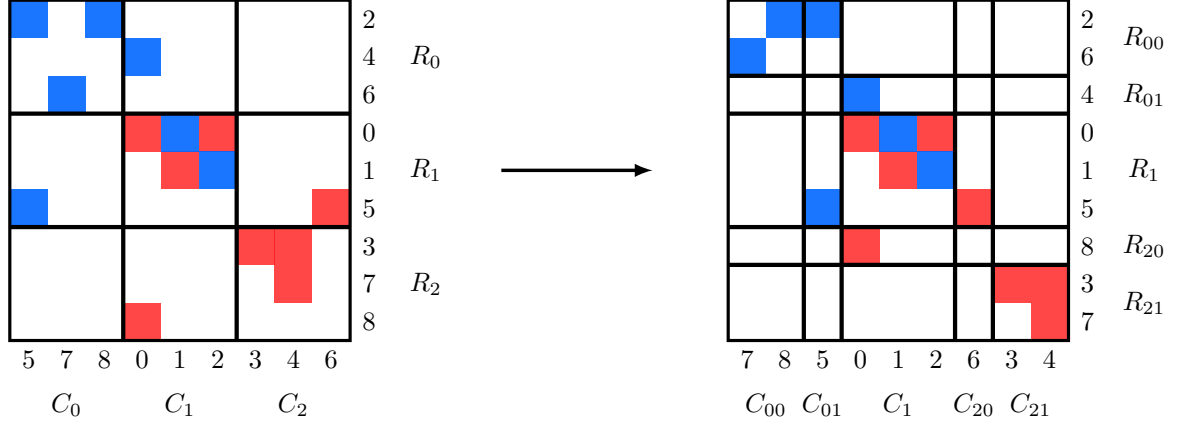


Figure 2.3: SBD2 form obtained starting from the SBD form of Figure 2.1.

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 full rows assigned to A_r and full 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 full assignments. The proposed heuristic does substantially this, by trying to swap the assignment of nonzeros from A_r to A_c and viceversa, trying to obtain that full rows are assigned to A_r and full columns are assigned to A_c . In order to obtain both of these things 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 fully assigning nonzeros in them to A_c) and of A_c , thus assigning full 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 little buffer: we pre-determine l_{max} , the maximum amount of worsening allowed, and, after this threshold is reached, we start considering only strictly improving solution. In order to have a meaningful threshold, it might be convenient to have it relative to the amount 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 the diagonal element (which correspond to not fully assigned rows/columns of A) with the minimum amount of nonzeros: our immediate goal, which in reality spans over a few moves of our local search, is to fully 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 2.6:

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: partitioned matrix A

Output: A_r, A_c

```

 $\ddot{A} :=$  SBD2 form of the partitioned matrix  $A$ 
 $m_i :=$  number of rows of  $\ddot{A}_{i*}$ 
 $n_i :=$  number of columns of  $\ddot{A}_{*i}$ 
 $A_r := A_c := \emptyset$ 
Assign nonzeros of  $\ddot{A}_{00}$  to  $A_r$ 
Assign nonzeros of  $\ddot{A}_{01}$  to  $A_r$ 
Assign nonzeros of  $\ddot{A}_{10}$  to  $A_c$ 
if  $m_1 < n_1$  then
    Assign nonzeros of  $\ddot{A}_{11}$  to  $A_c$ 
else
    Assign nonzeros of  $\ddot{A}_{11}$  to  $A_r$ 
end if
Assign nonzeros of  $\ddot{A}_{12}$  to  $A_r$ 
Assign nonzeros of  $\ddot{A}_{21}$  to  $A_c$ 
Assign nonzeros of  $\ddot{A}_{22}$  to  $A_r$  or  $A_c$  following Algorithm 2.1
Assign nonzeros of  $\ddot{A}_{23}$  to  $A_c$ 
Assign nonzeros of  $\ddot{A}_{32}$  to  $A_r$ 
if  $m_3 < n_3$  then
    Assign nonzeros of  $\ddot{A}_{33}$  to  $A_c$ 
else
    Assign nonzeros of  $\ddot{A}_{33}$  to  $A_r$ 
end if
Assign nonzeros of  $\ddot{A}_{34}$  to  $A_c$ 
Assign nonzeros of  $\ddot{A}_{43}$  to  $A_r$ 
Assign nonzeros of  $\ddot{A}_{44}$  to  $A_c$ 

```

Algorithm 2.5: Assignment of the nonzeros of A based on the SBD2 form of the matrix.

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 on 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 to be 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 , 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

Input: $A_r, A_c, l_{max}, iter_{max}$
Output: A'_r, A'_c

```

 $l := 0$ 
Compute  $B$  following the medium-grain model
for  $it = 1, \dots, iter_{max}$  do
   $i := \underset{k \in \{1, \dots, m+n\} \text{ s.t. } B(k,k) \neq 0}{\operatorname{argmin}} nz(k)$ 
  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
 $A'_r := B([1, \dots, n], [n+1, \dots, n+m])^T$ 
 $A'_c := B([n+1, \dots, m+n], [1, \dots, n])$ 

```

Algorithm 2.6: Local search refinement of A_r and A_c

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, we have that 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 2.7.

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 the row  $i$  to  $A_r$ 
  else
    Add the nonzeros of the column  $i$  to  $A_c$ 
  end if
end for

```

Algorithm 2.7: 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 2.8:

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

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

Algorithm 2.8: Alternative formulation of Algorithm 2.7.

Lastly, another different formulation is possible: we consider individually each nonzero a_{ij} and see whether in v we have $i < j$ (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.

A very important thing to observe is that in any case 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 scheme 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 divide 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 in 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.

In addition, we can refine a bit our sorting: we could move the indices which have only one nonzero to the back, because no matter our assignment of such nonzero, that index will not be cut and it is best to try and keep also the other dimension uncut.

- **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 we have twice as many columns than rows: in the first case we have

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

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 at 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.

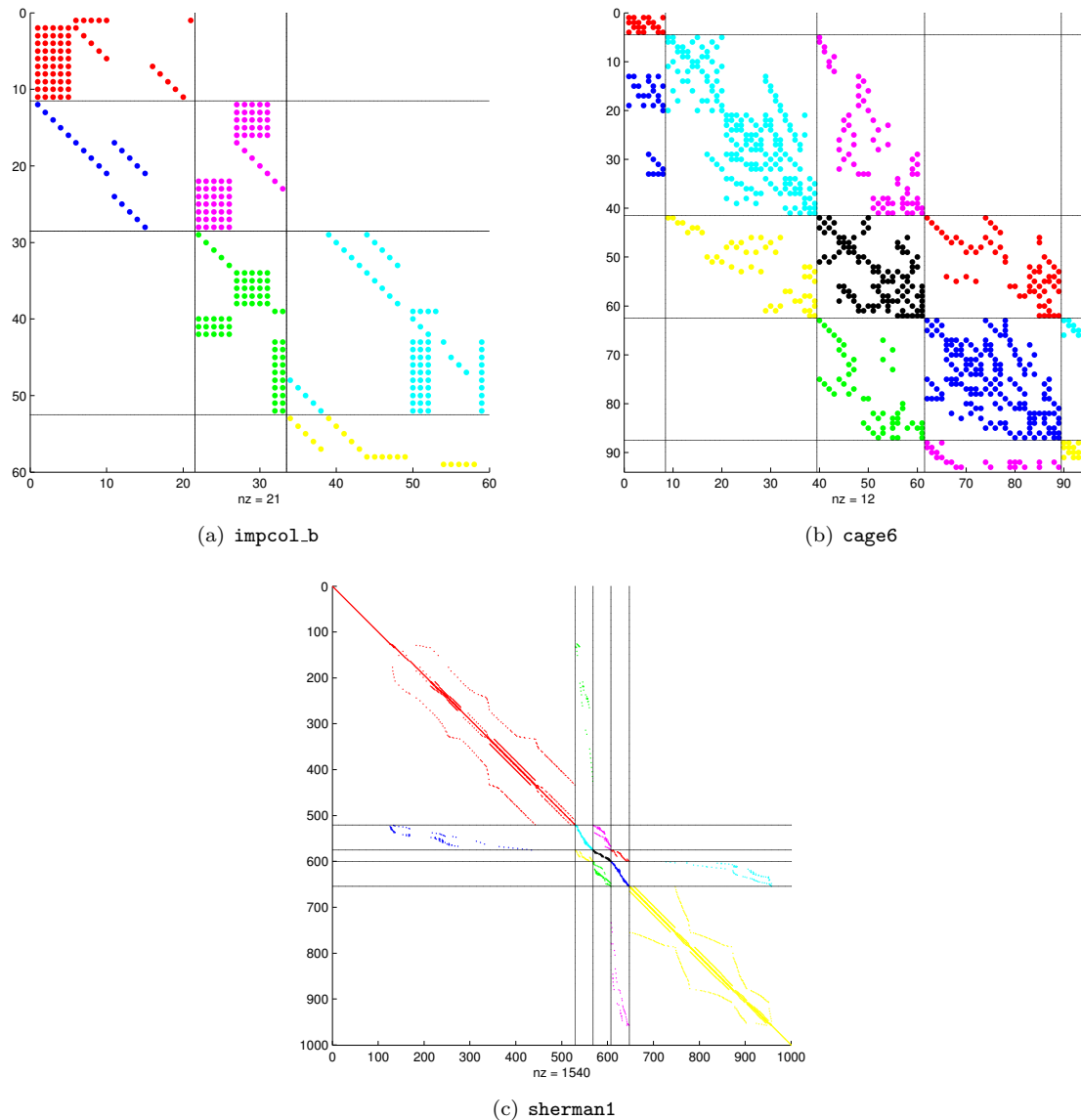


Figure 2.4: Example of SBD2 forms of three different matrices. Similarly as in Figure 2.2, each part of \tilde{A} has been given a color (note that since there are more parts than color used, some colors are repeated even though the parts are not related in any way). We can see in 2.4(a) that the second 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 DNA electrophoresis problem [25], while the latter is an oil reservoir simulation challenge matrix [24]. 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.

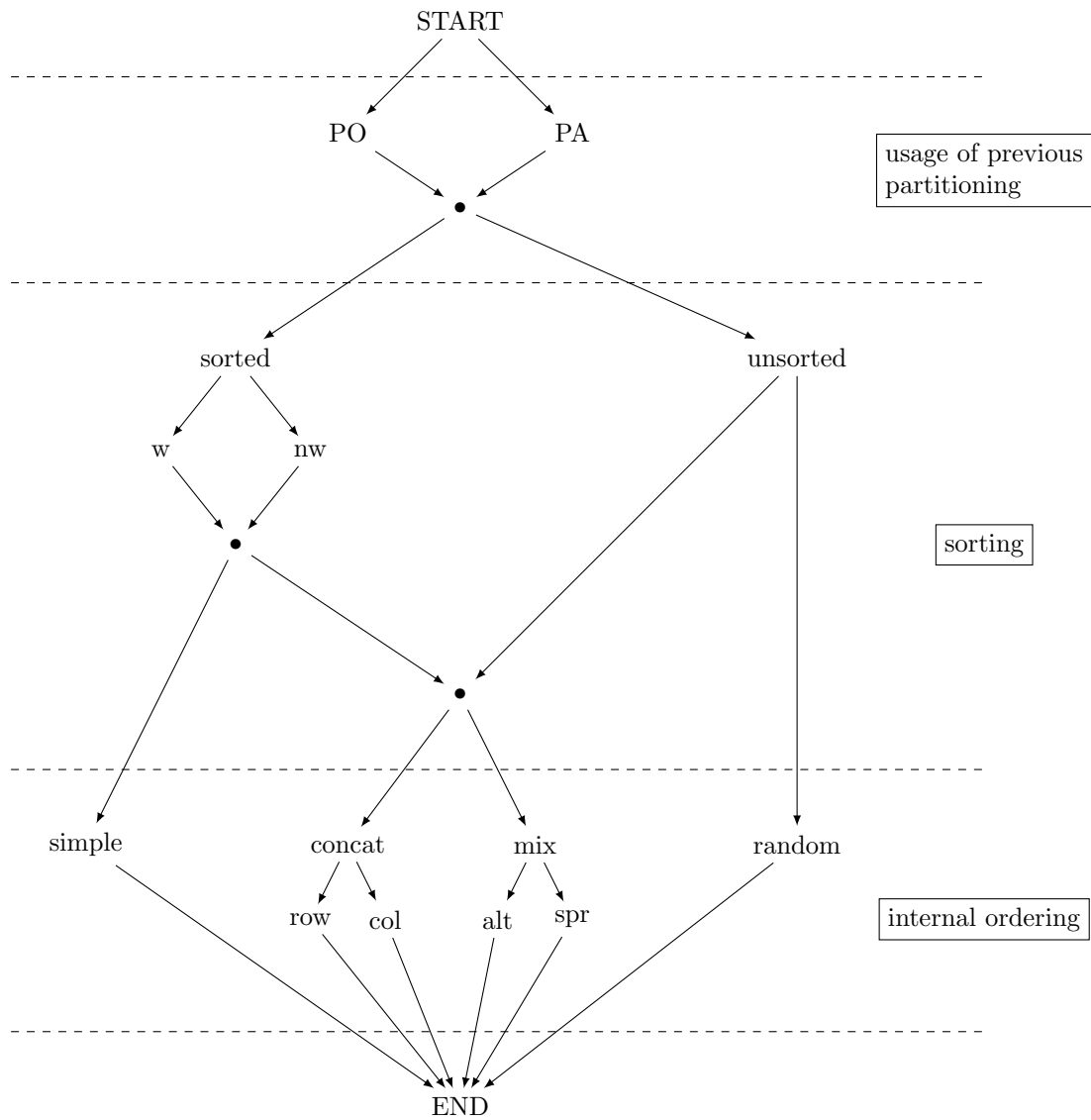


Figure 2.5: Directed graph that represents the family of heuristics used (any path from START to END). Dummy nodes (the ones without any label) were added in order to reduce the number of edges and ease legibility.

Chapter 3

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 lead 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, as already argued, the blocks \ddot{A}_{00} and \ddot{A}_{44} are interesting, because they contain nonzeros with the property of being “independent”. More precisely, those were rows and columns fully assigned to a processor and whose nonzeros didn’t have any neighbor (a nonzero in the same row or column) which had a cut column/row. These nonzeros, then, could be assigned anywhere that they still did not cause communication.

The term *independent*, in this reasoning, has to be defined very carefully: we want to look for a subset of $\{0, \dots, m+n-1\}$ which does not 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 much nonzeros as possible in this way, lowering consistently the upper bound on the communication volume that can be computed during the creation of A_r and A_c .

To do so, we can employ a very well studied object in graph theory: the **independent set**. However, this requires a correct translation of our sparse matrix into a graph, in order to produce the desired result, which is described in Section 3.1. In Section 3.2, instead, we discuss the actual algorithm used to compute the maximum independent set in such graph.

3.1 Graph construction

In order to retrieve the desired information from the graph, we need to construct it correctly from our sparse matrix. 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 the matrix partitioning problem, under the name of *bipartite graph model* by Hendrickson and Kolda [26], where the authors, after constructing the graph, discuss different algorithms for bipartite graph partitioning and come to the conclusion that the best strategy is using multilevel methods with Fiduccia-Mattheyses refinement.

More specifically, in this graph formulation we have that rows and columns are vertices, and we have the

edge (i, j) if $a_{ij} \neq 0$; the resulting graph is bipartite.

An example of such translation from matrix to graph is shown in Figure 3.1, where we translate the matrix given in Figure 1.1.

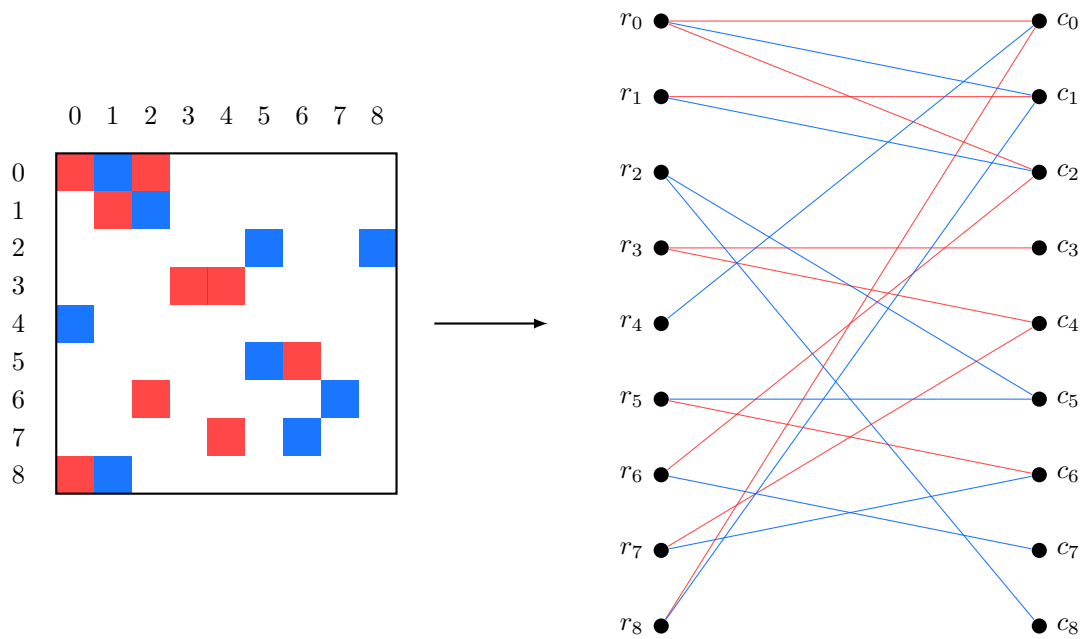


Figure 3.1: Graph constructed using the sparsity pattern of the matrix of Figure 1.1 as adjacency matrix (rows and columns are vertices, nonzeros are edges). The edge color is the same of the corresponding nonzeros, but only to facilitate the understanding of this translation; in reality there is no distinction between edges. With r_i we denote row i , whereas with c_j we denote column j .

3.2 Computation of the independent set

3.2.1 Hopcroft-Karp algorithm

Chapter 4

Experimental results

4.1 Improving the initial partitioning

4.2 Fully iterative partitioning

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).
- [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] Sivasankaran Rajamanickam and Erik G. Boman. “Parallel partitioning with Zoltan: Is hypergraph partitioning worth it?” In: *Graph Partitioning and Graph Clustering*. 2012, pp. 37–52 (cited on page 6).
- [7] Bruce Hendrickson. “Graph partitioning and parallel solvers: Has the emperor no clothes?” In: *Solving Irregularly Structured Problems in Parallel*. Springer, 1998, pp. 218–225 (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] G. Karypis, R. Aggarwal, V. Kumar, and Shashi Shekhar. “Multilevel hypergraph partitioning: applications in VLSI domain”, in: *Very Large Scale Integration (VLSI) Systems, IEEE Transactions on* 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: *in: Proceedings of the 15th International Parallel and Distributed Processing Symposium, 2001, p. 118. C. Aykanat*, pp. 609–625 (cited on pages 7, 8).
- [13] Daniël M. Pelt and Rob H. Bisseling. “A Medium-Grain Model for Fast 2D Bipartitioning of Sparse Matrices”, submitted. 2013 (cited on pages 8–10, 12).
- [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] Y.F. Hu, K.C.F. Maguire, and R.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] B.W. Kernighan and S. Lin. “An efficient heuristic procedure for partitioning graphs”, in: *Bell system technical journal* 49 (1970), pp. 291–307 (cited on page 9).
- [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 14, 16).
- [24] 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 16, 23).
- [25] 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 16, 23).
- [26] 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 25).