



# Optimizing nonzero-based sparse matrix partitioning models via reducing latency<sup>☆</sup>

Seher Acer, Oguz Selvitopi, Cevdet Aykanat<sup>\*</sup>

Bilkent University, Computer Engineering Department, 06800, Ankara, Turkey

## HIGHLIGHTS

- Optimizing fine-grain hypergraph model to reduce bandwidth and latency.
- Optimizing medium-grain hypergraph model to reduce bandwidth and latency.
- Message net concept to encapsulate minimization of total message count.
- Practical enhancements to establish a trade-off between bandwidth and latency.
- Significant performance improvements validated on nearly one thousand matrices.

## ARTICLE INFO

### Article history:

Received 22 November 2017  
Received in revised form 18 May 2018  
Accepted 5 August 2018  
Available online 18 August 2018

### Keywords:

Sparse matrix  
Sparse matrix–vector multiplication  
Row–column–parallel SpMV  
Load balancing  
Communication overhead  
Hypergraph  
Fine-grain partitioning  
Medium-grain partitioning  
Recursive bipartitioning

## ABSTRACT

For the parallelization of sparse matrix–vector multiplication (SpMV) on distributed memory systems, nonzero-based fine-grain and medium-grain partitioning models attain the lowest communication volume and computational imbalance among all partitioning models. This usually comes, however, at the expense of high message count, i.e., high latency overhead. This work addresses this shortcoming by proposing new fine-grain and medium-grain models that are able to minimize communication volume and message count in a single partitioning phase. The new models utilize message nets in order to encapsulate the minimization of total message count. We further fine-tune these models by proposing delayed addition and thresholding for message nets in order to establish a trade-off between the conflicting objectives of minimizing communication volume and message count. The experiments on an extensive dataset of nearly one thousand matrices show that the proposed models improve the total message count of the original nonzero-based models by up to 27% on the average, which is reflected on the parallel runtime of SpMV as an average reduction of 15% on 512 processors.

© 2018 Elsevier Inc. All rights reserved.

## 1. Introduction

Sparse matrix partitioning plays a pivotal role in scaling applications that involve irregular sparse matrices on distributed memory systems. Several decades of research on this subject led to elegant combinatorial partitioning models that are able to address the needs of these applications.

A key operation in sparse applications is the sparse matrix–vector multiplication (SpMV), which is usually performed in a repeated manner with the same sparse matrix in various iterative solvers. The irregular sparsity pattern of the matrix in SpMV necessitates a non-trivial parallelization. In that sense, graph and

hypergraph models prove to be powerful tools in their immense ability to represent SpMV with the aim of optimizing desired parallel performance metrics. We focus on the hypergraph models as they correctly encapsulate the total communication volume in SpMV [6,7,13,15] and the proposed models in this work rely on hypergraphs.

Among various hypergraph models, the fine-grain hypergraph model [8,10] achieves the lowest communication volume and the lowest imbalance on computational loads of the processors [10]. Since the nonzeros of the matrix are treated individually in the fine-grain model, the nonzeros that belong to the same row/column are more likely to be scattered to multiple processors compared to the other models. This may result in a high message count and hinder scalability. The fine-grain hypergraphs have the largest size for the same reason, causing this model to have the highest partitioning overhead. The recently proposed medium-grain model [19] alleviates this issue by operating on groups of nonzeros instead of individual nonzeros. The partitioning overhead

<sup>☆</sup> This work was supported by The Scientific and Technological Research Council of Turkey (TUBITAK) under Grant EEEAG-114E545. This article is also based upon work from COST Action CA 15109 (COSTNET).

<sup>\*</sup> Corresponding author.

E-mail addresses: [acer@cs.bilkent.edu.tr](mailto:acer@cs.bilkent.edu.tr) (S. Acer), [reha@cs.bilkent.edu.tr](mailto:reha@cs.bilkent.edu.tr) (O. Selvitopi), [aykanat@cs.bilkent.edu.tr](mailto:aykanat@cs.bilkent.edu.tr) (C. Aykanat).

of the medium-grain model is significantly lower than that of the fine-grain model, while these two models achieve comparable communication volume. The fine-grain and medium-grain models are referred to as nonzero-based models as they obtain nonzero-based matrix partitions, the most general possible [24].

Although the nonzero-based models attain the lowest communication volume, the overall communication cost is not determined by the volume only, but better formulated as a function of multiple communication cost metrics. Another important cost metric is the total message count, which is not only overlooked by both the fine-grain and medium-grain models, but also exacerbated due to having nonzero-based partitions. Note that among the two basic components of the communication cost, the total communication volume determines the bandwidth component and the total message count determines the latency component.

In this work, we aim at addressing the latency overheads of nonzero-based partitioning models. Our contributions can be summarized as follows:

- We propose a novel fine-grain model to simultaneously reduce the bandwidth and latency costs of parallel SpMV.
- We propose a novel medium-grain model to simultaneously reduce the bandwidth and latency costs of parallel SpMV.
- We utilize the message net concept [21] within the recursive bipartitioning framework to incorporate the minimization of the latency cost into the partitioning objective of these two models. Message nets aim to group the matrix nonzeros and/or the vector entries in the SpMV that necessitate a message together.
- We also propose two enhancements, delayed addition and thresholding for message nets, to better exploit the trade-off between the bandwidth and latency costs for the proposed models.
- We conduct extensive experiments on nearly one thousand matrices and show that the proposed models improve the total message count of the original nonzero-based models by up to 27% on the average, which is reflected on the parallel runtime of SpMV as an average reduction of 15% on 512 processors.

The remainder of the paper is organized as follows. Section 2 gives background on parallel SpMV, performance cost metrics, the fine-grain model, recursive bipartitioning, and the medium-grain model. Sections 3 and 4 present the proposed fine-grain and medium-grain models, respectively. Section 5 describes practical enhancements to these models. Sections 6 and 7 give the experimental results and related work, respectively, and Section 8 concludes.

## 2. Preliminaries

### 2.1. Row-column-parallel SpMV

We consider the parallelization of SpMV of the form  $y = Ax$  with a nonzero-based partitioned matrix  $A$ , where  $A = (a_{ij})$  is an  $n_r \times n_c$  sparse matrix with  $n_{nz}$  nonzero entries, and  $x$  and  $y$  are dense vectors. The  $i$ th row and the  $j$ th column of  $A$  are, respectively, denoted by  $r_i$  and  $c_j$ . The  $j$ th entry of  $x$  and the  $i$ th entry of  $y$  are, respectively, denoted by  $x_j$  and  $y_i$ . Let  $\mathcal{A}$  denote the set of nonzero entries in  $A$ , that is,  $\mathcal{A} = \{a_{ij} : a_{ij} \neq 0\}$ . Let  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively, denote the sets of entries in  $x$  and  $y$ , that is,  $\mathcal{X} = \{x_1, \dots, x_{n_c}\}$  and  $\mathcal{Y} = \{y_1, \dots, y_{n_r}\}$ . Assume that there are  $K$  processors in the parallel system denoted by  $P_1, \dots, P_K$ . Let  $\Pi_K(\mathcal{A}) = \{\mathcal{A}_1, \dots, \mathcal{A}_K\}$ ,  $\Pi_K(\mathcal{X}) = \{\mathcal{X}_1, \dots, \mathcal{X}_K\}$ , and  $\Pi_K(\mathcal{Y}) = \{\mathcal{Y}_1, \dots, \mathcal{Y}_K\}$  denote  $K$ -way partitions of  $\mathcal{A}$ ,  $\mathcal{X}$ , and  $\mathcal{Y}$ , respectively.

Given partitions  $\Pi_K(\mathcal{A})$ ,  $\Pi_K(\mathcal{X})$ , and  $\Pi_K(\mathcal{Y})$ , without loss of generality, the nonzeros in  $\mathcal{A}_k$  and the vector entries in  $\mathcal{X}_k$  and  $\mathcal{Y}_k$  are

**Algorithm 1:** Row-column-parallel SpMV as performed by processor  $P_k$ .

**Require:**  $\mathcal{A}_k, \mathcal{X}_k$

---

▷ *Pre-communication phase* — expands on  $x$ -vector entries  
 Receive the needed  $x$ -vector entries that are not in  $\mathcal{X}_k$   
 Send the  $x$ -vector entries in  $\mathcal{X}_k$  needed by other processors

▷ *Computation phase*  
 $y_i^{(k)} \leftarrow y_i^{(k)} + a_{ij}x_j$  for each  $a_{ij} \in \mathcal{A}_k$

▷ *Post-communication phase* — folds on  $y$ -vector entries  
 Receive the partial results for  $y$ -vector entries in  $\mathcal{Y}_k$  and  
 compute  $y_i \leftarrow \sum y_i^{(\ell)}$  for each partial result  $y_i^{(\ell)}$   
 Send the partial results for  $y$ -vector entries not in  $\mathcal{Y}_k$

**return**  $\mathcal{Y}_k$

---

assigned to processor  $P_k$ . For each  $a_{ij} \in \mathcal{A}_k$ ,  $P_k$  is held responsible for performing the respective multiply-and-add operation  $y_i^{(k)} \leftarrow y_i^{(k)} + a_{ij}x_j$ , where  $y_i^{(k)}$  denotes the partial result computed for  $y_i$  by  $P_k$ . Algorithm 1 displays the basic steps performed by  $P_k$  in parallel SpMV for a nonzero-based partitioned matrix  $A$ . This algorithm is called the *row-column-parallel SpMV* [22]. In this algorithm,  $P_k$  first receives the needed  $x$ -vector entries that are not in  $\mathcal{X}_k$  from their owners and sends its  $x$ -vector entries to the processors that need them in a *pre-communication* phase. Sending  $x_j$  to possibly multiple processors is referred to as the *expand* operation on  $x_j$ . When  $P_k$  has all needed  $x$ -vector entries, it performs the local SpMV by computing  $y_i^{(k)} \leftarrow y_i^{(k)} + a_{ij}x_j$  for each  $a_{ij} \in \mathcal{A}_k$ .  $P_k$  then receives the partial results for the  $y$ -vector entries in  $\mathcal{Y}_k$  from other processors and sends its partial results to the processors that own the respective  $y$ -vector entries in a *post-communication* phase. Receiving partial result(s) for  $y_i$  from possibly multiple processors is referred to as the *fold* operation on  $y_i$ . Note that overlapping of computation and communication is not considered in this algorithm for the sake of clarity.

### 2.2. Performance cost metrics

In this section, we describe the performance cost metrics that are minimized by the proposed models and formulate them on given  $K$ -way partitions  $\Pi_K(\mathcal{A})$ ,  $\Pi_K(\mathcal{X})$ , and  $\Pi_K(\mathcal{Y})$ . Let  $e(P_k, P_\ell)$  denote the set of  $x$ -vector entries sent (expanded) from processor  $P_k$  to processor  $P_\ell$  during the pre-communication phase. Similarly, let  $f(P_k, P_\ell)$  denote the set of partial results for  $y$ -vector entries sent (folded) from  $P_k$  to  $P_\ell$  during the post-communication phase. That is,

$$\begin{aligned} e(P_k, P_\ell) &= \{x_j : x_j \in \mathcal{X}_k \text{ and } \exists a_{t,j} \in \mathcal{A}_\ell\} \text{ and} \\ f(P_k, P_\ell) &= \{y_i^{(k)} : y_i \in \mathcal{Y}_\ell \text{ and } \exists a_{i,t} \in \mathcal{A}_k\}. \end{aligned} \quad (1)$$

*Total communication volume* is equal to the sum of the sizes of all messages transmitted during pre-communication and post-communication phases and formulated as

$$\sum_k \sum_\ell |e(P_k, P_\ell)| + |f(P_k, P_\ell)|.$$

*Total message count* is equal to the total number of messages and formulated as

$$|\{(k, \ell) : e(P_k, P_\ell) \neq \emptyset\}| + |\{(k, \ell) : f(P_k, P_\ell) \neq \emptyset\}|.$$

*Computational imbalance* is equal to the ratio of the maximum to the average amount of computation performed by a processor minus one. Since the amount of computation in SpMV is proportional to the number of nonzeros, computational imbalance is formulated

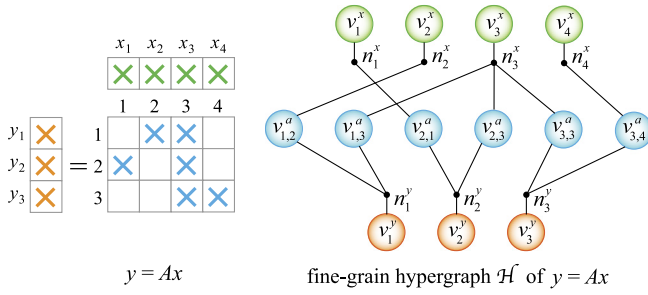


Fig. 1. A sample  $y = Ax$  and the corresponding fine-grain hypergraph.

as

$$\frac{\max_k |\mathcal{A}_k|}{|\mathcal{A}|/K} - 1.$$

For an efficient row-column-parallel SpMV, the goal is to find partitions  $\Pi_K(\mathcal{A})$ ,  $\Pi_K(\mathcal{X})$  and  $\Pi_K(\mathcal{Y})$  that achieve low communication overheads and low computational imbalance. Existing fine-grain [9] and medium-grain [19] models, which are, respectively, described in Sections 2.3 and 2.5, meet this goal partially by only minimizing the bandwidth cost (i.e., total communication volume) while maintaining balance on the computational loads of processors.

### 2.3. Fine-grain hypergraph model

In the fine-grain hypergraph  $\mathcal{H} = (\mathcal{V}, \mathcal{N})$ , each entry in  $\mathcal{A}$ ,  $\mathcal{X}$ , and  $\mathcal{Y}$  is represented by a different vertex. Vertex set  $\mathcal{V}$  contains a vertex  $v_{i,j}^a$  for each  $a_{i,j} \in \mathcal{A}$ , a vertex  $v_j^x$  for each  $x_j \in \mathcal{X}$ , and a vertex  $v_i^y$  for each  $y_i \in \mathcal{Y}$ . That is,

$$\mathcal{V} = \{v_{i,j}^a : a_{i,j} \neq 0\} \cup \{v_1^x, \dots, v_{n_c}^x\} \cup \{v_1^y, \dots, v_{n_r}^y\}.$$

$v_{i,j}^a$  represents both the data element  $a_{i,j}$  and the computational task  $y_i \leftarrow y_i + a_{i,j}x_j$  associated with  $a_{i,j}$ , whereas  $v_j^x$  and  $v_i^y$  only represent the input and output data elements  $x_j$  and  $y_i$ , respectively.

The net set  $\mathcal{N}$  contains two different types of nets to represent the dependencies of the computational tasks on  $x$ - and  $y$ -vector entries. For each  $x_j \in \mathcal{X}$  and  $y_i \in \mathcal{Y}$ ,  $\mathcal{N}$ , respectively, contains the nets  $n_j^x$  and  $n_i^y$ . That is,

$$\mathcal{N} = \{n_1^x, \dots, n_{n_c}^x\} \cup \{n_1^y, \dots, n_{n_r}^y\}.$$

Net  $n_j^x$  represents the input dependency of the computational tasks on  $x_j$ ; hence, it connects the vertices that represent these tasks and  $v_j^x$ . Net  $n_i^y$  represents the output dependency of the computational tasks on  $y_i$ ; hence, it connects the vertices that represent these tasks and  $v_i^y$ . The sets of vertices connected by  $n_j^x$  and  $n_i^y$  are, respectively, formulated as

$$\text{Pins}(n_j^x) = \{v_{i,j}^a\} \cup \{v_{i,t}^a : a_{i,t} \neq 0\} \text{ and}$$

$$\text{Pins}(n_i^y) = \{v_i^y\} \cup \{v_{i,t}^a : a_{i,t} \neq 0\}.$$

$\mathcal{H}$  contains  $n_{nz} + n_c + n_r$  vertices,  $n_c + n_r$  nets and  $2n_{nz} + n_c + n_r$  pins. Fig. 1 displays a sample SpMV instance and its corresponding fine-grain hypergraph. In  $\mathcal{H}$ , the vertices are assigned the weights that signify their computational loads. Hence,  $w(v_{i,j}^a) = 1$  for each  $v_{i,j}^a \in \mathcal{V}$  as  $v_{i,j}^a$  represents a single multiply-and-add operation, whereas  $w(v_j^x) = w(v_i^y) = 0$  for each  $v_j^x \in \mathcal{V}$  and  $v_i^y \in \mathcal{V}$  as they do not represent any computation. The nets are assigned unit costs, i.e.,  $c(n_j^x) = c(n_i^y) = 1$  for each  $n_j^x \in \mathcal{N}$  and  $n_i^y \in \mathcal{N}$ .

A  $K$ -way vertex partition  $\Pi_K(\mathcal{H}) = \{\mathcal{V}_1, \dots, \mathcal{V}_K\}$  can be decoded to obtain  $\Pi_K(\mathcal{A})$ ,  $\Pi_K(\mathcal{X})$ , and  $\Pi_K(\mathcal{Y})$  by assigning the entries

represented by the vertices in part  $\mathcal{V}_k$  to processor  $P_k$ . That is,

$$\mathcal{A}_k = \{a_{i,j} : v_{i,j}^a \in \mathcal{V}_k\},$$

$$\mathcal{X}_k = \{x_j : v_j^x \in \mathcal{V}_k\}, \text{ and}$$

$$\mathcal{Y}_k = \{y_i : v_i^y \in \mathcal{V}_k\}.$$

Let  $\Lambda(n)$  denote the set of the parts connected by net  $n$  in  $\Pi_K(\mathcal{H})$ , where a net is said to connect a part if it connects at least one vertex in that part. Let  $\lambda(n)$  denote the number of parts connected by  $n$ , i.e.,  $|\Lambda(n)|$ . A net  $n$  is called cut if it connects at least two parts, i.e.,  $\lambda(n) > 1$ , and uncut, otherwise. The cutsize of  $\Pi_K(\mathcal{H})$  is defined as

$$\text{cutsizesize}(\Pi_K(\mathcal{H})) = \sum_{n \in \mathcal{N}} c(n)(\lambda(n) - 1). \quad (2)$$

Consider cut nets  $n_j^x$  and  $n_i^y$  in  $\Pi_K(\mathcal{H})$  and assume that  $v_j^x, v_i^y \in \mathcal{V}_k$ . The cut net  $n_j^x$  necessitates sending (expanding) of  $x_j$  from  $P_k$  to the processors that correspond to the parts in  $\Lambda(n_j^x) - \{\mathcal{V}_k\}$  in the pre-communication phase. Hence, it can be said that if  $n_j^x$  is cut, then  $v_j^x$  incurs a communication volume of  $\lambda(n_j^x) - 1$ . The cut net  $n_i^y$ , on the other hand, necessitates sending (folding) of the partial results from the processors that correspond to the parts in  $\Lambda(n_i^y) - \{\mathcal{V}_k\}$  to  $P_k$ , which then sums them up to obtain  $y_i$ . Hence, it can be said that if  $n_i^y$  is cut, then  $v_i^y$  incurs a communication volume of  $\lambda(n_i^y) - 1$ . Since each cut net  $n$  increases the cutsize by  $\lambda(n) - 1 > 0$ ,  $\text{cutsizesize}(\Pi_K(\mathcal{H}))$  is equal to the sum of the volume in pre- and post-communication phases. Therefore, minimizing  $\text{cutsizesize}(\Pi_K(\mathcal{H}))$  corresponds to minimizing the total communication volume in row-column-parallel SpMV.

The contents of the messages sent from  $P_k$  to  $P_\ell$  in the pre- and post-communication phases in terms of a partitioned hypergraph are, respectively, given by

$$e(P_k, P_\ell) = \{x_j : v_j^x \in \mathcal{V}_k \text{ and } v_\ell \in \Lambda(n_j^x)\} \text{ and} \quad (3)$$

$$f(P_k, P_\ell) = \{y_i^{(k)} : v_i^y \in \mathcal{V}_\ell \text{ and } v_k \in \Lambda(n_i^y)\}.$$

Note the one-to-one correspondence between sparse matrix and hypergraph partitionings in determining the message contents given by Eqs. (1) and (3).

In  $\Pi_K(\mathcal{H})$ , the weight  $W(\mathcal{V}_k)$  of part  $\mathcal{V}_k$  is defined as the sum of the weights of the vertices in  $\mathcal{V}_k$ , i.e.,  $W(\mathcal{V}_k) = \sum_{v \in \mathcal{V}_k} w(v)$ , which is equal to the total computational load of processor  $P_k$ . Then, maintaining the balance constraint

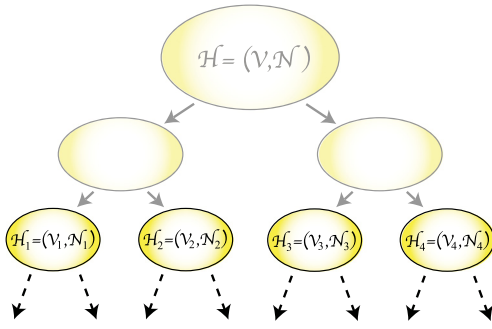
$$W(\mathcal{V}_k) \leq W_{\text{avg}}(1 + \epsilon), \text{ for } k = 1, \dots, K,$$

corresponds to maintaining balance on the computational loads of the processors. Here,  $W_{\text{avg}}$  and  $\epsilon$  denote the average part weight and a maximum imbalance ratio, respectively.

### 2.4. Recursive bipartitioning (RB) paradigm

In RB, a given domain is first bipartitioned and then this bipartition is used to form two new subdomains. In our case, a domain refers to a hypergraph ( $\mathcal{H}$ ) or a set of matrix and vector entries ( $\mathcal{A}, \mathcal{X}, \mathcal{Y}$ ). The newly-formed subdomains are recursively bipartitioned until  $K$  subdomains are obtained. This procedure forms a hypothetical full binary tree, which contains  $\lceil \log K \rceil + 1$  levels. The root node of the tree represents the given domain, whereas each of the remaining nodes represents a subdomain formed during the RB process. At any stage of the RB process, the subdomains represented by the leaf nodes of the RB tree collectively induce a partition of the original domain.

The RB paradigm is successfully used for hypergraph partitioning. Fig. 2 illustrates an RB tree currently in the process of partitioning a hypergraph. The current leaf nodes induce a four-way partition  $\Pi_4(\mathcal{H}) = \{\mathcal{V}_1, \mathcal{V}_2, \mathcal{V}_3, \mathcal{V}_4\}$  and each node in the RB tree represents both a hypergraph and its vertex set. While forming



**Fig. 2.** The RB tree during partitioning  $\mathcal{H} = (\mathcal{V}, \mathcal{N})$ . The current RB tree contains four leaf hypergraphs with the hypergraph to be bipartitioned next being  $\mathcal{H}_1 = (\mathcal{V}_1, \mathcal{N}_1)$ .

two new subhypergraphs after each RB step, the cut-net splitting technique is used [7] to encapsulate the cutsizes in (2). The sum of the cutsizes incurred in all RB steps is equal to the cutsize of the resulting  $K$ -way partition.

### 2.5. Medium-grain hypergraph model

In the medium-grain hypergraph model, the sets  $\mathcal{A}$ ,  $\mathcal{X}$  and  $\mathcal{Y}$  are partitioned into  $K$  parts using RB. The medium-grain model uses a mapping for a subset of the nonzeros at each RB step. Because this mapping is central to the model, we focus on a single bipartitioning step to explain the medium-grain model. Before each RB step, the nonzeros to be bipartitioned are first mapped to their rows or columns by a heuristic and a new hypergraph is formed according to this mapping.

Consider an RB tree for the medium-grain model with  $K'$  leaf nodes, where  $K' < K$ , and assume that the  $k$ th node from the left is to be bipartitioned next. This node represents  $\mathcal{A}_k$ ,  $\mathcal{X}_k$ , and  $\mathcal{Y}_k$  in the respective  $K'$ -way partitions  $\{\mathcal{A}_1, \dots, \mathcal{A}_{K'}\}$ ,  $\{\mathcal{X}_1, \dots, \mathcal{X}_{K'}\}$ , and  $\{\mathcal{Y}_1, \dots, \mathcal{Y}_{K'}\}$ . First, each  $a_{i,j} \in \mathcal{A}_k$  is mapped to either  $r_i$  or  $c_j$ , where this mapping is denoted by  $\text{map}(a_{i,j})$ . With a heuristic,  $a_{i,j} \in \mathcal{A}_k$  is mapped to  $r_i$  if  $r_i$  has fewer nonzeros than  $c_j$  in  $\mathcal{A}_k$ , and to  $c_j$  if  $c_j$  has fewer nonzeros than  $r_i$  in  $\mathcal{A}_k$ . After determining  $\text{map}(a_{i,j})$  for each nonzero in  $\mathcal{A}_k$ , the medium-grain hypergraph  $\mathcal{H}_k = (\mathcal{V}_k, \mathcal{N}_k)$  is formed as follows. Vertex set  $\mathcal{V}_k$  contains a vertex  $v_j^x$  if  $x_j$  is in  $\mathcal{X}_k$  or there exists at least one nonzero in  $\mathcal{A}_k$  mapped to  $c_j$ . Similarly,  $\mathcal{V}_k$  contains a vertex  $v_i^y$  if  $y_i$  is in  $\mathcal{Y}_k$  or there exists at least one nonzero in  $\mathcal{A}_k$  mapped to  $r_i$ . Hence,  $v_j^x$  represents  $x_j$  and/or the nonzero(s) assigned to  $c_j$ , whereas  $v_i^y$  represents  $y_i$  and/or the nonzero(s) assigned to  $r_i$ . That is,

$$\mathcal{V}_k = \{v_j^x : x_j \in \mathcal{X}_k \text{ or } \exists a_{t,j} \in \mathcal{A}_k \text{ s.t. } \text{map}(a_{t,j}) = c_j\} \cup \{v_i^y : y_i \in \mathcal{Y}_k \text{ or } \exists a_{i,t} \in \mathcal{A}_k \text{ s.t. } \text{map}(a_{i,t}) = r_i\}.$$

Besides the data elements, vertex  $v_j^x/v_i^y$  represents the group of computational tasks associated with the nonzeros mapped to them, if any.

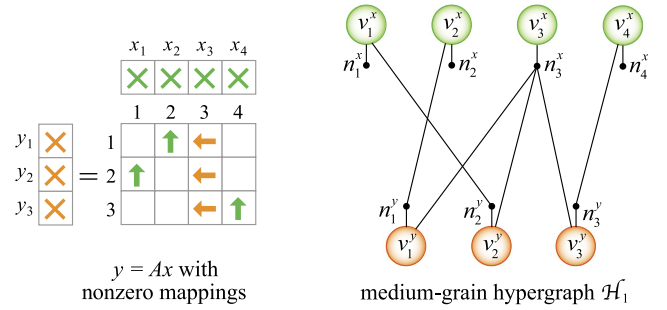
The net set  $\mathcal{N}_k$  contains a net  $n_j^x$  if  $\mathcal{A}_k$  contains at least one nonzero in  $c_j$ , and a net  $n_i^y$  if  $\mathcal{A}_k$  contains at least one nonzero in  $r_i$ . That is,

$$\mathcal{N}_k = \{n_j^x : \exists a_{t,j} \in \mathcal{A}_k\} \cup \{n_i^y : \exists a_{i,t} \in \mathcal{A}_k\}.$$

$n_j^x$  represents the input dependency of the groups of computational tasks on  $x_j$ , whereas  $n_i^y$  represents the output dependency of the groups of computational tasks on  $y_i$ . Hence, the sets of vertices connected by  $n_j^x$  and  $n_i^y$  are, respectively, formulated by

$$\text{Pins}(n_j^x) = \{v_j^x\} \cup \{v_t^y : \text{map}(a_{t,j}) = r_t\} \text{ and}$$

$$\text{Pins}(n_i^y) = \{v_i^y\} \cup \{v_t^x : \text{map}(a_{i,t}) = c_t\}.$$



**Fig. 3.** The nonzero assignments of the sample  $y = Ax$  and the corresponding medium-grain hypergraph.

In  $\mathcal{H}_k$ , each net is assigned a unit cost, i.e.,  $c(n_j^x) = c(n_i^y) = 1$  for each  $n_j^x \in \mathcal{N}$  and  $n_i^y \in \mathcal{N}$ . Each vertex is assigned a weight equal to the number of nonzeros represented by that vertex. That is,

$$w(v_j^x) = |\{a_{t,j} : \text{map}(a_{t,j}) = c_j\}| \text{ and}$$

$$w(v_i^y) = |\{a_{i,t} : \text{map}(a_{i,t}) = r_i\}|.$$

$\mathcal{H}_k$  is bipartitioned with the objective of minimizing the cutsize and the constraint of maintaining balance on the part weights. The resulting bipartition is further improved by an iterative refinement algorithm. In every RB step, minimizing the cutsize corresponds to minimizing the total volume of communication, whereas maintaining balance on the weights of the parts corresponds to maintaining balance on the computational loads of the processors.

Fig. 3 displays a sample SpMV instance with nonzero mapping information and the corresponding medium-grain hypergraph. This example illustrates the first RB step, hence,  $\mathcal{A}_1 = \mathcal{A}$ ,  $\mathcal{X}_1 = \mathcal{X}$ ,  $\mathcal{Y}_1 = \mathcal{Y}$ , and  $K' = k = 1$ . Each nonzero in  $A$  is denoted by an arrow, where the direction of the arrow shows the mapping for that nonzero. For example,  $n_3^x$  connects  $v_3^x$ ,  $v_1^y$ ,  $v_2^y$ , and  $v_3^y$  since  $\text{map}(a_{1,3}) = r_1$ ,  $\text{map}(a_{2,3}) = r_2$ , and  $\text{map}(a_{3,3}) = r_3$ .

### 3. Optimizing fine-grain partitioning model

In this section, we propose a fine-grain hypergraph partitioning model that simultaneously reduces the bandwidth and latency costs of the row-column-parallel SpMV. Our model is built upon the original fine-grain model (Section 2.3) via utilizing the RB paradigm. The proposed model contains two different types of nets to address the bandwidth and latency costs. The nets of the original fine-grain model already address the bandwidth cost and they are called “volume nets” as they encapsulate the minimization of the total communication volume. At each RB step, our model forms and adds new nets to the hypergraph to be bipartitioned. These new nets address the latency cost and they are called “message nets” as they encapsulate the minimization of the total message count.

Message nets aim to group the matrix nonzeros and vector entries that altogether necessitate a message. The formation and addition of message nets rely on the RB paradigm. To determine the existence and the content of a message, a partition information is needed first. At each RB step, prior to bipartitioning the current hypergraph that already contains the volume nets, the message nets are formed using the  $K'$ -way partition information and added to this hypergraph, where  $K'$  is the number of leaf nodes in the current RB tree. Then this hypergraph is bipartitioned, which results in a  $(K' + 1)$ -way partition as the number of leaves becomes  $K' + 1$  after bipartitioning. Adding message nets just before each bipartitioning allows us to utilize the most recent global partition information at hand. In contrast to the formation of the message nets, the formation of the volume nets via cut-net splitting requires only the local bipartition information.



### 3.1. Message nets in a single RB step

Consider an SpMV instance  $y = Ax$  and its corresponding fine-grain hypergraph  $\mathcal{H} = (\mathcal{V}, \mathcal{N})$  with the aim of partitioning  $\mathcal{H}$  into  $K$  parts to parallelize  $y = Ax$ . The RB process starts with bipartitioning  $\mathcal{H}$ , which is represented by the root node of the corresponding RB tree. Assume that the RB process is at the state where there are  $K'$  leaf nodes in the RB tree, for  $1 < K' < K$ , and the hypergraphs corresponding to these nodes are denoted by  $\mathcal{H}_1, \dots, \mathcal{H}_{K'}$  from left to right. Let  $\Pi_{K'}(\mathcal{H}) = \{\mathcal{V}_1, \dots, \mathcal{V}_{K'}\}$  denote the  $K'$ -way partition induced by the leaf nodes of the RB tree.  $\Pi_{K'}(\mathcal{H})$  also induces  $K'$ -way partitions  $\Pi_{K'}(\mathcal{A})$ ,  $\Pi_{K'}(\mathcal{X})$ , and  $\Pi_{K'}(\mathcal{Y})$  of sets  $\mathcal{A}$ ,  $\mathcal{X}$ , and  $\mathcal{Y}$ , respectively. Without loss of generality, the entries in  $\mathcal{A}_k$ ,  $\mathcal{X}_k$ , and  $\mathcal{Y}_k$  are assigned to processor group  $\mathcal{P}_k$ . Assume that  $\mathcal{H}_k = (\mathcal{V}_k, \mathcal{N}_k)$  is next to be bipartitioned among these hypergraphs.  $\mathcal{H}_k$  initially contains only the volume nets. In our model, we add message nets to  $\mathcal{H}_k$  to obtain the augmented hypergraph  $\mathcal{H}_k^M = (\mathcal{V}_k, \mathcal{N}_k^M)$ . Let  $\Pi(\mathcal{H}_k^M) = \{\mathcal{V}_{k,L}, \mathcal{V}_{k,R}\}$  denote a bipartition of  $\mathcal{H}_k^M$ , where  $L$  and  $R$  in the subscripts refer to left and right, respectively.  $\Pi(\mathcal{H}_k^M)$  induces bipartitions  $\Pi(\mathcal{A}_k) = \{\mathcal{A}_{k,L}, \mathcal{A}_{k,R}\}$ ,  $\Pi(\mathcal{X}_k) = \{\mathcal{X}_{k,L}, \mathcal{X}_{k,R}\}$ , and  $\Pi(\mathcal{Y}_k) = \{\mathcal{Y}_{k,L}, \mathcal{Y}_{k,R}\}$  on  $\mathcal{A}_k$ ,  $\mathcal{X}_k$ , and  $\mathcal{Y}_k$ , respectively. Let  $\mathcal{P}_{k,L}$  and  $\mathcal{P}_{k,R}$  denote the processor groups to which the entries in  $\{\mathcal{A}_{k,L}, \mathcal{X}_{k,L}, \mathcal{Y}_{k,L}\}$  and  $\{\mathcal{A}_{k,R}, \mathcal{X}_{k,R}, \mathcal{Y}_{k,R}\}$  are assigned.

Algorithm 2 displays the basic steps of forming message nets and adding them to  $\mathcal{H}_k$ . For each processor group  $\mathcal{P}_\ell$  that  $\mathcal{P}_k$  communicates with, four different message nets may be added to  $\mathcal{H}_k$ : expand-send net, expand-receive net, fold-send net and fold-receive net, respectively, denoted by  $s_\ell^e$ ,  $r_\ell^e$ ,  $s_\ell^f$  and  $r_\ell^f$ . Here,  $s$  and  $r$ , respectively, denote the messages sent and received, the subscript  $\ell$  denotes the id of the processor group communicated with, and the superscripts  $e$  and  $f$ , respectively, denote the expand and fold operations. These nets are next explained in detail.

- **expand-send net  $s_\ell^e$** : Net  $s_\ell^e$  represents the message sent from  $\mathcal{P}_k$  to  $\mathcal{P}_\ell$  during the expand operations on  $x$ -vector entries in the pre-communication phase. This message consists of the  $x$ -vector entries owned by  $\mathcal{P}_k$  and needed by  $\mathcal{P}_\ell$ . Hence,  $s_\ell^e$  connects the vertices that represent the  $x$ -vector entries required by the computational tasks in  $\mathcal{P}_\ell$ . That is,

$$\text{Pins}(s_\ell^e) = \{v_j^x : x_j \in \mathcal{X}_k \text{ and } \exists a_{t,j} \in \mathcal{A}_\ell\}.$$

The formation and addition of expand-send nets are performed in lines 2–7 of Algorithm 2. After bipartitioning  $\mathcal{H}_k^M$ , if  $s_\ell^e$  becomes cut in  $\Pi(\mathcal{H}_k^M)$ , both  $\mathcal{P}_{k,L}$  and  $\mathcal{P}_{k,R}$  send a message to  $\mathcal{P}_\ell$ , where the contents of the messages sent from  $\mathcal{P}_{k,L}$  and  $\mathcal{P}_{k,R}$  to  $\mathcal{P}_\ell$  are  $\{x_j : v_j^x \in \mathcal{V}_{k,L} \text{ and } a_{t,j} \in \mathcal{A}_\ell\}$  and  $\{x_j : v_j^x \in \mathcal{V}_{k,R} \text{ and } a_{t,j} \in \mathcal{A}_\ell\}$ , respectively. The overall number of messages in the pre-communication phase increases by one in this case since  $\mathcal{P}_k$  was sending a single message to  $\mathcal{P}_\ell$  and it is split into two messages after bipartitioning. If  $s_\ell^e$  becomes uncut, the overall number of messages does not change since only one of  $\mathcal{P}_{k,L}$  and  $\mathcal{P}_{k,R}$  sends a message to  $\mathcal{P}_\ell$ .

- **expand-receive net  $r_\ell^e$** : Net  $r_\ell^e$  represents the message received by  $\mathcal{P}_k$  from  $\mathcal{P}_\ell$  during the expand operations on  $x$ -vector entries in the pre-communication phase. This message consists of the  $x$ -vector entries owned by  $\mathcal{P}_\ell$  and needed by  $\mathcal{P}_k$ . Hence,  $r_\ell^e$  connects the vertices that represent the computational tasks requiring  $x$ -vector entries from  $\mathcal{P}_\ell$ . That is,

$$\text{Pins}(r_\ell^e) = \{v_{t,j}^a : a_{t,j} \in \mathcal{A}_k \text{ and } x_j \in \mathcal{X}_\ell\}.$$

The formation and addition of expand-receive nets are performed in lines 8–13 of Algorithm 2. After bipartitioning  $\mathcal{H}_k^M$ , if  $r_\ell^e$  becomes cut in  $\Pi(\mathcal{H}_k^M)$ , both  $\mathcal{P}_{k,L}$  and  $\mathcal{P}_{k,R}$  receive a message from  $\mathcal{P}_\ell$ , where the contents of the messages received by  $\mathcal{P}_{k,L}$  and  $\mathcal{P}_{k,R}$  from  $\mathcal{P}_\ell$  are  $\{x_j : v_{t,j}^a \in \mathcal{V}_{k,L} \text{ and } x_j \in \mathcal{X}_\ell\}$  and  $\{x_j : v_{t,j}^a \in \mathcal{V}_{k,R} \text{ and } x_j \in \mathcal{X}_\ell\}$ , respectively. The overall number

#### Algorithm 2: ADD-MESSAGE-NETS.

---

**Require:**  $\mathcal{H}_k = (\mathcal{V}_k, \mathcal{N}_k)$ ,  $\Pi_{K'}(\mathcal{A}) = \{\mathcal{A}_1, \dots, \mathcal{A}_{K'}\}$ ,  $\Pi_{K'}(\mathcal{X}) = \{\mathcal{X}_1, \dots, \mathcal{X}_{K'}\}$ ,  $\Pi_{K'}(\mathcal{Y}) = \{\mathcal{Y}_1, \dots, \mathcal{Y}_{K'}\}$ .

- 1:  $\mathcal{N}_k^M \leftarrow \mathcal{N}_k$
- $\triangleright$  Expand-send nets
- 2: **for** each  $x_j \in \mathcal{X}_k$  **do**
- 3:   **for** each  $a_{t,j} \in \mathcal{A}_{\ell \neq k}$  **do**
- 4:     **if**  $s_\ell^e \notin \mathcal{N}_k^M$  **then**
- 5:        $\text{Pins}(s_\ell^e) \leftarrow \{v_j^x\}$ ,  $\mathcal{N}_k^M \leftarrow \mathcal{N}_k^M \cup \{s_\ell^e\}$
- 6:     **else**
- 7:        $\text{Pins}(s_\ell^e) \leftarrow \text{Pins}(s_\ell^e) \cup \{v_j^x\}$
- $\triangleright$  Expand-receive nets
- 8: **for** each  $a_{t,j} \in \mathcal{A}_k$  **do**
- 9:   **for** each  $x_j \in \mathcal{X}_{\ell \neq k}$  **do**
- 10:     **if**  $r_\ell^e \notin \mathcal{N}_k^M$  **then**
- 11:        $\text{Pins}(r_\ell^e) \leftarrow \{v_{t,j}^a\}$ ,  $\mathcal{N}_k^M \leftarrow \mathcal{N}_k^M \cup \{r_\ell^e\}$
- 12:     **else**
- 13:        $\text{Pins}(r_\ell^e) \leftarrow \text{Pins}(r_\ell^e) \cup \{v_{t,j}^a\}$
- $\triangleright$  Fold-send nets
- 14: **for** each  $a_{i,t} \in \mathcal{A}_k$  **do**
- 15:   **for** each  $y_i \in \mathcal{Y}_{\ell \neq k}$  **do**
- 16:     **if**  $s_\ell^f \notin \mathcal{N}_k^M$  **then**
- 17:        $\text{Pins}(s_\ell^f) \leftarrow \{v_{i,t}^a\}$ ,  $\mathcal{N}_k^M \leftarrow \mathcal{N}_k^M \cup \{s_\ell^f\}$
- 18:     **else**
- 19:        $\text{Pins}(s_\ell^f) \leftarrow \text{Pins}(s_\ell^f) \cup \{v_{i,t}^a\}$
- $\triangleright$  Fold-receive nets
- 20: **for** each  $y_i \in \mathcal{Y}_k$  **do**
- 21:   **for** each  $a_{i,t} \in \mathcal{A}_{\ell \neq k}$  **do**
- 22:     **if**  $r_\ell^f \notin \mathcal{N}_k^M$  **then**
- 23:        $\text{Pins}(r_\ell^f) \leftarrow \{v_i^y\}$ ,  $\mathcal{N}_k^M \leftarrow \mathcal{N}_k^M \cup \{r_\ell^f\}$
- 24:     **else**
- 25:        $\text{Pins}(r_\ell^f) \leftarrow \text{Pins}(r_\ell^f) \cup \{v_i^y\}$
- 26: **return**  $\mathcal{H}_k^M = (\mathcal{V}_k, \mathcal{N}_k^M)$

---

of messages in the pre-communication phase increases by one in this case and does not change if  $r_\ell^e$  becomes uncut.

- **fold-send net  $s_\ell^f$** : Net  $s_\ell^f$  represents the message sent from  $\mathcal{P}_k$  to  $\mathcal{P}_\ell$  during the fold operations on  $y$ -vector entries in the post-communication phase. This message consists of the partial results computed by  $\mathcal{P}_k$  for the  $y$ -vector entries owned by  $\mathcal{P}_\ell$ . Hence,  $s_\ell^f$  connects the vertices that represent the computational tasks whose partial results are required by  $\mathcal{P}_\ell$ . That is,

$$\text{Pins}(s_\ell^f) = \{v_{i,t}^a : a_{i,t} \in \mathcal{A}_k \text{ and } y_i \in \mathcal{Y}_\ell\}.$$

The formation and addition of fold-send nets are performed in lines 14–19 of Algorithm 2. After bipartitioning  $\mathcal{H}_k^M$ , if  $s_\ell^f$  becomes cut in  $\Pi(\mathcal{H}_k^M)$ , both  $\mathcal{P}_{k,L}$  and  $\mathcal{P}_{k,R}$  send a message to  $\mathcal{P}_\ell$ , where the contents of the messages sent from  $\mathcal{P}_{k,L}$  and  $\mathcal{P}_{k,R}$  to  $\mathcal{P}_\ell$  are  $\{y_i^{(k,L)} : v_{i,t}^a \in \mathcal{V}_{k,L} \text{ and } y_i \in \mathcal{Y}_\ell\}$  and  $\{y_i^{(k,R)} : v_{i,t}^a \in \mathcal{V}_{k,R} \text{ and } y_i \in \mathcal{Y}_\ell\}$ , respectively. The overall number of messages in the post-communication phase increases by one in this case and does not change if  $s_\ell^f$  becomes uncut.

- **fold-receive net  $r_\ell^f$** : Net  $r_\ell^f$  represents the message received by  $\mathcal{P}_k$  from  $\mathcal{P}_\ell$  during the fold operations on  $y$ -vector entries in the post-communication phase. This message consists of the partial results computed by  $\mathcal{P}_\ell$  for the  $y$ -vector entries owned by  $\mathcal{P}_k$ . Hence,  $r_\ell^f$  connects the vertices that represent the  $y$ -vector entries for which  $\mathcal{P}_\ell$  produces partial results. That is,

$$\text{Pins}(r_\ell^f) = \{v_i^y : y_i \in \mathcal{Y}_k \text{ and } \exists a_{i,t} \in \mathcal{A}_\ell\}.$$

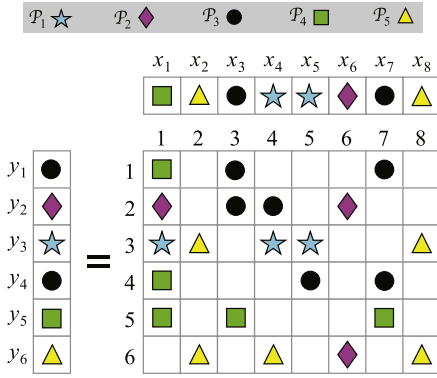


Fig. 4. A 5-way nonzero-based partition of an SpMV instance  $y = Ax$ .

Table 1

The messages communicated by  $\mathcal{P}_3$  in pre- and post-communication phases before and after bipartitioning  $\mathcal{H}_3^M$ . The number of messages communicated by  $\mathcal{P}_3$  increases from 4 to 6 due to two cut message nets in  $\Pi(\mathcal{H}_3^M)$ .

RB state	Phase	Message	Due to
Before $\Pi(\mathcal{H}_3^M)$	Pre	$\mathcal{P}_3$ sends $\{x_3, x_7\}$ to $\mathcal{P}_4$	$a_{5,3}, a_{5,7}$
		$\mathcal{P}_3$ receives $\{x_4, x_5\}$ from $\mathcal{P}_1$	$a_{2,4}, a_{4,5}$
	Post	$\mathcal{P}_3$ sends $\{y_2^{(3)}\}$ to $\mathcal{P}_2$	$a_{2,3}, a_{2,4}$
		$\mathcal{P}_3$ receives $\{y_1^{(4)}, y_4^{(4)}\}$ from $\mathcal{P}_4$	$a_{1,1}, a_{4,1}$
After $\Pi(\mathcal{H}_3^M)$	Pre	$\mathcal{P}_{3,L}$ sends $\{x_3, x_7\}$ to $\mathcal{P}_4$	$a_{5,3}, a_{5,7}$
		$\mathcal{P}_{3,R}$ receives $\{x_4, x_5\}$ from $\mathcal{P}_1$	$a_{2,4}, a_{4,5}$
		$\mathcal{P}_{3,L}$ sends $\{y_2^{(3,L)}\}$ to $\mathcal{P}_2$	$a_{2,3}$
	Post	$\mathcal{P}_{3,R}$ sends $\{y_2^{(3,R)}\}$ to $\mathcal{P}_2$	$a_{2,4}$
		$\mathcal{P}_{3,L}$ receives $\{y_1^{(4)}\}$ from $\mathcal{P}_4$	$a_{1,1}$
		$\mathcal{P}_{3,R}$ receives $\{y_4^{(4)}\}$ from $\mathcal{P}_4$	$a_{4,1}$

The formation and addition of fold-receive nets are performed in lines 20–25 of Algorithm 2. After bipartitioning  $\mathcal{H}_k^M$ , if  $r_\ell^f$  becomes cut in  $\Pi(\mathcal{H}_k^M)$ , both  $\mathcal{P}_{k,L}$  and  $\mathcal{P}_{k,R}$  receive a message from  $\mathcal{P}_\ell$ , where the contents of the messages received by  $\mathcal{P}_{k,L}$  and  $\mathcal{P}_{k,R}$  from  $\mathcal{P}_\ell$  are  $\{y_i^{(\ell)} : v_i^y \in \mathcal{V}_{k,L} \text{ and } a_{i,\ell} \in \mathcal{A}_\ell\}$  and  $\{y_i^{(\ell)} : v_i^y \in \mathcal{V}_{k,R} \text{ and } a_{i,\ell} \in \mathcal{A}_\ell\}$ , respectively. The overall number of messages in the post-communication phase increases by one in this case and does not change if  $r_\ell^f$  becomes uncut.

Note that at most four message nets are required to encapsulate the messages between processor groups  $\mathcal{P}_k$  and  $\mathcal{P}_\ell$ . The message nets in  $\mathcal{H}_k^M$  encapsulate all the messages that  $\mathcal{P}_k$  communicates with other processor groups. Since the number of leaf hypergraphs is  $K'$ ,  $\mathcal{P}_k$  may communicate with at most  $K' - 1$  processor groups, hence the maximum number of message nets that can be added to  $\mathcal{H}_k$  is  $4(K' - 1)$ .

Fig. 4 displays an SpMV instance with a  $6 \times 8$  matrix  $A$ , which is being partitioned by the proposed model. The RB process is at the state where there are five leaf hypergraphs  $\mathcal{H}_1, \dots, \mathcal{H}_5$ , and the hypergraph to be bipartitioned next is  $\mathcal{H}_3$ . The figure displays the assignments of the matrix nonzeros and vector entries to the corresponding processor groups  $\mathcal{P}_1, \dots, \mathcal{P}_5$ . Each symbol in the figure represents a distinct processor group and a symbol inside a cell signifies the assignment of the corresponding matrix nonzero or vector entry to the processor group represented by that symbol. For example, the nonzeros in  $\mathcal{A}_3 = \{a_{1,3}, a_{1,7}, a_{2,3}, a_{2,4}, a_{4,5}, a_{4,7}\}$ ,  $x$ -vector entries in  $\mathcal{X}_3 = \{x_3, x_7\}$ , and  $y$ -vector entries in  $\mathcal{Y}_3 = \{y_1, y_4\}$  are assigned to  $\mathcal{P}_3$ . The left of Fig. 5 displays the augmented hypergraph  $\mathcal{H}_3^M$  that contains volume and message nets. In the figure, the volume nets are illustrated by small black circles with thin lines, whereas the message nets are illustrated by the respective processor's symbol with thick lines.

The messages communicated by  $\mathcal{P}_3$  under the assignments given in Fig. 4 are displayed at the top half of Table 1. In the pre-communication phase,  $\mathcal{P}_3$  sends a message to  $\mathcal{P}_4$  and receives a message from  $\mathcal{P}_1$ , and in the post-communication phase, it sends a message to  $\mathcal{P}_2$  and receives a message from  $\mathcal{P}_4$ . Hence, we add four message nets to  $\mathcal{H}_3$ : expand-send net  $s_4^e$ , expand-receive net  $r_1^e$ , fold-send net  $s_2^f$ , and fold-receive net  $r_4^f$ . In Fig. 5, for example,  $r_1^e$  connects the vertices  $v_{2,4}^a$  and  $v_{4,5}^a$  since it represents the message received by  $\mathcal{P}_3$  from  $\mathcal{P}_1$  containing  $\{x_4, x_5\}$  due to nonzeros  $a_{2,4}$  and  $a_{4,5}$ . The right of Fig. 5 displays a bipartition  $\Pi(\mathcal{H}_3^M)$  and the messages that  $\mathcal{P}_{3,L}$  and  $\mathcal{P}_{3,R}$  communicate with the other processor groups due to  $\Pi(\mathcal{H}_3^M)$  are given in the bottom half of Table 1. Since  $s_4^e$  and  $r_1^e$  are uncut, only one of  $\mathcal{P}_{3,L}$  and  $\mathcal{P}_{3,R}$  participates in sending or receiving the corresponding message. Since  $s_2^f$  is cut, both  $\mathcal{P}_{3,L}$  and  $\mathcal{P}_{3,R}$  send a message to  $\mathcal{P}_2$ , and since  $r_4^f$  is cut, both  $\mathcal{P}_{3,L}$  and  $\mathcal{P}_{3,R}$  receive a message from  $\mathcal{P}_4$ .

In  $\mathcal{H}_k^M$ , each volume net is assigned the cost of the per-word transfer time,  $t_w$ , whereas each message net is assigned the cost of the start-up latency,  $t_{su}$ . Let  $v$  and  $m$ , respectively, denote the number of volume and message nets that are cut in  $\Pi(\mathcal{H}_k^M)$ . Then,

$$\text{cutsizes}(\Pi(\mathcal{H}_k^M)) = vt_w + mt_{su}.$$

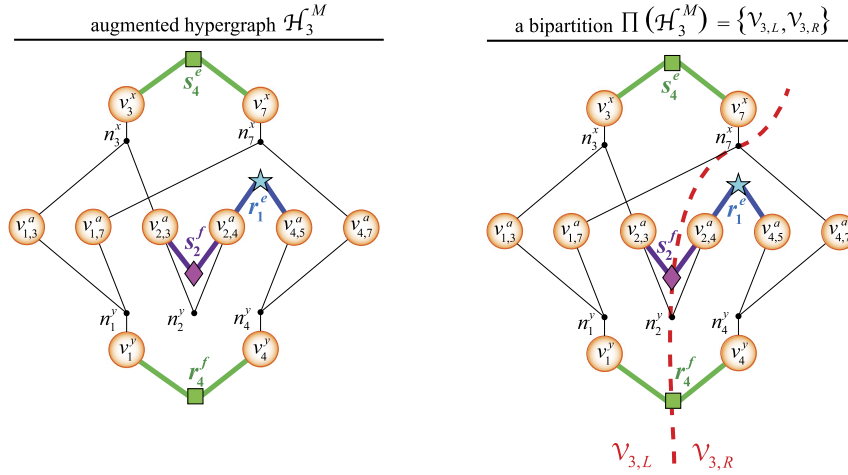
Here,  $v$  is equal to the increase in the total communication volume incurred by  $\Pi(\mathcal{H}_k^M)$  [7]. Recall that each cut message net increases the number of messages that  $\mathcal{P}_k$  communicates with the respective processor group by one. Hence,  $m$  is equal to the increase in the number of messages that  $\mathcal{P}_k$  communicates with other processor groups. The overall increase in the total message count due to  $\Pi(\mathcal{H}_k^M)$  is  $m + \delta$ , where  $\delta$  denotes the number of messages between  $\mathcal{P}_{k,L}$  and  $\mathcal{P}_{k,R}$ , and is bounded by two (empirically found to be almost always two). Hence, minimizing the cutsizes of  $\Pi(\mathcal{H}_k^M)$  corresponds to simultaneously reducing the increase in the total communication volume and the total message count in the respective RB step. Therefore, minimizing the cutsizes in all RB steps corresponds to reducing the total communication volume and the total message count simultaneously.

After obtaining a bipartition  $\Pi(\mathcal{H}_k^M) = \{\mathcal{V}_{k,L}, \mathcal{V}_{k,R}\}$  of the augmented hypergraph  $\mathcal{H}_k^M$ , the new hypergraphs  $\mathcal{H}_{k,L} = (\mathcal{V}_{k,L}, \mathcal{N}_{k,L})$  and  $\mathcal{H}_{k,R} = (\mathcal{V}_{k,R}, \mathcal{N}_{k,R})$  are immediately formed with only volume nets. Recall that the formation of the volume nets of  $\mathcal{H}_{k,L}$  and  $\mathcal{H}_{k,R}$  is performed with the cut-net splitting technique and it can be performed using the local bipartition information  $\Pi(\mathcal{H}_k^M)$ .

### 3.2. The overall RB

After completing an RB step and obtaining  $\mathcal{H}_{k,L}$  and  $\mathcal{H}_{k,R}$ , the labels of the hypergraphs represented by the leaf nodes of the RB tree are updated as follows. For  $1 \leq i < k$ , the label of  $\mathcal{H}_i = (\mathcal{V}_i, \mathcal{N}_i)$  does not change. For  $k < i < K'$ ,  $\mathcal{H}_i = (\mathcal{V}_i, \mathcal{N}_i)$  becomes  $\mathcal{H}_{i+1} = (\mathcal{V}_{i+1}, \mathcal{N}_{i+1})$ . Hypergraphs  $\mathcal{H}_{k,L} = (\mathcal{V}_{k,L}, \mathcal{N}_{k,L})$  and  $\mathcal{H}_{k,R} = (\mathcal{V}_{k,R}, \mathcal{N}_{k,R})$  become  $\mathcal{H}_k = (\mathcal{V}_k, \mathcal{N}_k)$  and  $\mathcal{H}_{k+1} = (\mathcal{V}_{k+1}, \mathcal{N}_{k+1})$ , respectively. As a result, the vertex sets corresponding to the updated leaf nodes induce a  $(K' + 1)$ -way partition  $\Pi_{K'+1}(\mathcal{H}) = \{\mathcal{V}_1, \dots, \mathcal{V}_{K'+1}\}$ . The RB process then continues with the next hypergraph  $\mathcal{H}_{k+2}$  to be bipartitioned, which was labeled with  $\mathcal{H}_{k+1}$  in the previous RB state.

We next provide the cost of adding message nets through Algorithm 2 in the entire RB process. For the addition of expand-send nets, all nonzeros  $a_{t,j} \in \mathcal{A}_{\ell \neq k}$  with  $x_j \in \mathcal{X}_k$  are visited once (lines 2–7). Since  $\mathcal{X}_k \cap \mathcal{X}_\ell = \emptyset$  for  $1 \leq k \neq \ell \leq K'$  and  $\mathcal{X} = \bigcup_{k=1}^{K'} \mathcal{X}_k$ , each nonzero of  $A$  is visited once. For the addition of expand-receive nets, all nonzeros in  $\mathcal{A}_k$  are visited once (lines 8–13). Hence, each nonzero of  $A$  is visited once during the bipartitionings in a level of the RB tree since  $\mathcal{A}_k \cap \mathcal{A}_\ell = \emptyset$  for  $1 \leq k \neq \ell \leq K'$  and  $\mathcal{A} = \bigcup_{k=1}^{K'} \mathcal{A}_k$ . Therefore, the cost of adding expand-send and



**Fig. 5.** Left: Augmented hypergraph  $\mathcal{H}_3^M$  with 5 volume and 4 message nets. Right: A bipartition  $\Pi(\mathcal{H}_3^M) = \{V_{3,L}, V_{3,R}\}$  with two cut message nets ( $s_2^f, r_4^f$ ) and two cut volume nets ( $n_7^x, n_2^y$ ).

expand-receive nets is  $O(n_{nz})$  in a single level of the RB tree. A dual discussion holds for the addition of fold-send and fold-receive nets. Since the RB tree contains  $\lceil \log K \rceil$  levels in which bipartitionings take place, the overall cost of adding message nets is  $O(n_{nz} \log K)$ .

### 3.3. Adaptation for conformal partitioning

Partitions on input and output vectors  $x$  and  $y$  are said to be conformal if  $x_i$  and  $y_i$  are assigned to the same processor, for  $1 \leq i \leq n_r = n_c$ . Note that conformal vector partitions are valid for  $y = Ax$  with a square matrix. The motivation for a conformal partition arises in iterative solvers in which the  $y_i$  in an iteration is used to compute the  $x_i$  of the next iteration via linear vector operations. Assigning  $x_i$  and  $y_i$  to the same processor prevents the redundant communication of  $y_i$  to the processor that owns  $x_i$ .

Our model does not impose conformal partitions on vectors  $x$  and  $y$ , i.e.,  $x_i$  and  $y_i$  can be assigned to different processors. However, it is possible to adapt our model to obtain conformal partitions on  $x$  and  $y$  using the vertex amalgamation technique proposed in [23]. To assign  $x_i$  and  $y_i$  to the same processor, the vertices  $v_i^x$  and  $v_i^y$  are amalgamated into a new vertex  $v_i^{x/y}$ , which represents both  $x_i$  and  $y_i$ . The weight of  $v_i^{x/y}$  is set to be zero since the weights of  $v_i^x$  and  $v_i^y$  are zero. In  $\mathcal{H}_k^M$ , each volume/message net that connects  $v_i^x$  or  $v_i^y$  now connects the amalgamated vertex  $v_i^{x/y}$ . At each RB step,  $x_i$  and  $y_i$  are both assigned to the processor group corresponding to the leaf hypergraph that contains  $v_i^{x/y}$ .

## 4. Optimizing medium-grain partitioning model

In this section, we propose a medium-grain hypergraph partitioning model that simultaneously reduces the bandwidth and latency costs of the row-column-parallel SpMV. Our model is built upon the original medium-grain partitioning model (Section 2.5). The medium-grain hypergraphs in RB are augmented with the message nets before they are bipartitioned as in the fine-grain model proposed in Section 3. Since the fine-grain and medium-grain models both obtain nonzero-based partitions, the types and meanings of the message nets used in the medium-grain model are the same as those used in the fine-grain model. However, forming message nets for a medium-grain hypergraph is more involved due to the mappings used in this model.

Consider an SpMV instance  $y = Ax$  and the corresponding sets  $\mathcal{A}$ ,  $\mathcal{X}$ , and  $\mathcal{Y}$ . Assume that the RB process is at the state before bipartitioning the  $k$ th leaf node where there are  $K'$  leaf nodes in the current RB tree. Recall from Section 2.5 that the leaf nodes induce  $K'$ -way partitions  $\Pi_{K'}(\mathcal{A}) = \{\mathcal{A}_1, \dots, \mathcal{A}_{K'}\}$ ,  $\Pi_{K'}(\mathcal{X}) =$

$\{\mathcal{X}_1, \dots, \mathcal{X}_{K'}\}$  and  $\Pi_{K'}(\mathcal{Y}) = \{\mathcal{Y}_1, \dots, \mathcal{Y}_{K'}\}$ , and the  $k$ th leaf node represents  $\mathcal{A}_k$ ,  $\mathcal{X}_k$ , and  $\mathcal{Y}_k$ . To obtain bipartitions of  $\mathcal{A}_k$ ,  $\mathcal{X}_k$ , and  $\mathcal{Y}_k$ , we perform the following four steps.

(1) *Form the medium-grain hypergraph  $\mathcal{H}_k = (\mathcal{V}_k, \mathcal{N}_k)$  using  $\mathcal{A}_k$ ,  $\mathcal{X}_k$ , and  $\mathcal{Y}_k$ .* This process is the same with that in the original medium-grain model (Section 2.5). Recall that the nets in the medium-grain hypergraph encapsulate the total communication volume. Hence, these nets are assigned a cost of  $t_w$ .

(2) *Add message nets to  $\mathcal{H}_k$  to obtain augmented hypergraph  $\mathcal{H}_k^M$ .* For each processor group  $\mathcal{P}_\ell$  other than  $\mathcal{P}_k$ , there are four possible message nets that can be added to  $\mathcal{H}_k$ :

- **expand-send net  $s_\ell^e$ :** The set of vertices connected by  $s_\ell^e$  is the same with that of the expand-send net in the fine-grain model.
- **expand-receive net  $r_\ell^e$ :** The set of vertices connected by  $r_\ell^e$  is given by
 
$$\text{Pins}(r_\ell^e) = \{v_j^x : \exists a_{t,j} \in \mathcal{A}_k \text{ s.t. } \text{map}(a_{t,j}) = c_j \text{ and } x_j \in \mathcal{X}_\ell\} \cup \{v_i^y : \exists a_{t,i} \in \mathcal{A}_k \text{ s.t. } \text{map}(a_{t,i}) = r_i \text{ and } y_i \in \mathcal{X}_\ell\}.$$
- **fold-send net  $s_\ell^f$ :** The set of vertices connected by  $s_\ell^f$  is given by
 
$$\text{Pins}(s_\ell^f) = \{v_t^x : \exists a_{i,t} \in \mathcal{A}_k \text{ s.t. } \text{map}(a_{i,t}) = c_t \text{ and } y_i \in \mathcal{Y}_\ell\} \cup \{v_i^y : \exists a_{i,t} \in \mathcal{A}_k \text{ s.t. } \text{map}(a_{i,t}) = r_i \text{ and } y_i \in \mathcal{Y}_\ell\}.$$
- **fold-receive net  $r_\ell^f$ :** The set of vertices connected by  $r_\ell^f$  is the same with that of the fold-receive net in the fine-grain model.

The message nets are assigned a cost of  $t_{su}$  as they encapsulate the latency cost.

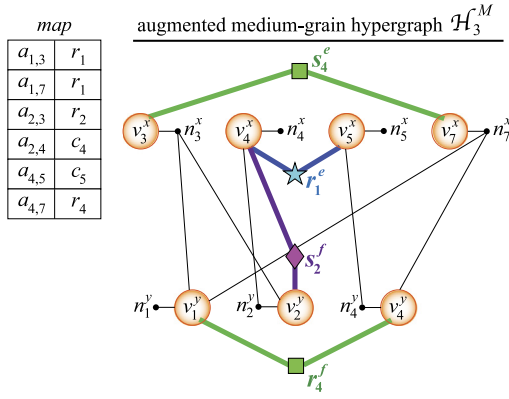
(3) *Obtain a bipartition  $\Pi(\mathcal{H}_k^M)$ .*  $\mathcal{H}_k^M$  is bipartitioned to obtain  $\Pi(\mathcal{H}_k^M) = \{V_{k,L}, V_{k,R}\}$ .

(4) *Derive bipartitions  $\Pi(\mathcal{A}_k) = \{\mathcal{A}_{k,L}, \mathcal{A}_{k,R}\}$ ,  $\Pi(\mathcal{X}_k) = \{\mathcal{X}_{k,L}, \mathcal{X}_{k,R}\}$  and  $\Pi(\mathcal{Y}_k) = \{\mathcal{Y}_{k,L}, \mathcal{Y}_{k,R}\}$  from  $\Pi(\mathcal{H}_k^M)$ .* For each nonzero  $a_{i,j} \in \mathcal{A}_k$ ,  $a_{i,j}$  is assigned to  $\mathcal{A}_{k,L}$  if the vertex that represents  $a_{i,j}$  is in  $V_{k,L}$ , and to  $\mathcal{A}_{k,R}$ , otherwise. That is,

$$\begin{aligned} \mathcal{A}_{k,L} &= \{a_{i,j} : \text{map}(a_{i,j}) = c_j \text{ with } v_j^x \in V_{k,L} \text{ or} \\ &\quad \text{map}(a_{i,j}) = r_i \text{ with } v_i^y \in V_{k,L}\} \text{ and} \\ \mathcal{A}_{k,R} &= \{a_{i,j} : \text{map}(a_{i,j}) = c_j \text{ with } v_j^x \in V_{k,R} \text{ or} \\ &\quad \text{map}(a_{i,j}) = r_i \text{ with } v_i^y \in V_{k,R}\}. \end{aligned}$$

For each  $x$ -vector entry  $x_j \in \mathcal{X}_k$ ,  $x_j$  is assigned to  $\mathcal{X}_{k,L}$  if  $v_j^x \in V_{k,L}$ , and to  $\mathcal{X}_{k,R}$ , otherwise. That is,

$$\mathcal{X}_{k,L} = \{x_j : v_j^x \in V_{k,L}\} \text{ and } \mathcal{X}_{k,R} = \{x_j : v_j^x \in V_{k,R}\}.$$



**Fig. 6.** The augmented medium-grain hypergraph  $\mathcal{H}_3^M$  formed during the RB process for the SpMV instance given in Fig. 4.

Similarly, for each  $y$ -vector entry  $y_i \in \mathcal{Y}_k$ ,  $y_i$  is assigned to  $\mathcal{Y}_{k,L}$  if  $v_i^y \in \mathcal{V}_{k,L}$ , and to  $\mathcal{Y}_{k,R}$ , otherwise. That is,

$$\mathcal{Y}_{k,L} = \{y_i : v_i^y \in \mathcal{V}_{k,L}\} \text{ and } \mathcal{Y}_{k,R} = \{y_i : v_i^y \in \mathcal{V}_{k,R}\}.$$

Fig. 6 displays the medium-grain hypergraph  $\mathcal{H}_3^M = (\mathcal{V}_3, \mathcal{N}_3^M)$  augmented with message nets, which is formed during bipartitioning  $\mathcal{A}_3$ ,  $\mathcal{X}_3$  and  $\mathcal{Y}_3$  given in Fig. 4. The table in the figure displays  $\text{map}(a_{i,j})$  value for each nonzero in  $\mathcal{A}_3$  computed by the heuristic described in Section 2.5. Augmented medium-grain hypergraph  $\mathcal{H}_3^M$  has four message nets. Observe that the sets of vertices connected by expand-send net  $s_4^e$  and fold-receive net  $r_4^f$  are the same for the fine-grain and medium-grain hypergraphs, which are, respectively, illustrated in Figs. 5 and 6. Expand-receive net  $r_1^e$  connects  $v_4^x$  and  $v_5^x$  since  $\mathcal{P}_3$  receives  $\{x_4, x_5\}$  due to nonzeros in  $\{a_{2,4}, a_{4,5}\}$  with  $\text{map}(a_{2,4}) = c_4$  and  $\text{map}(a_{4,5}) = c_5$ . Fold-send net  $s_2^f$  connects  $v_4^x$  and  $v_2^y$  since  $\mathcal{P}_3$  sends partial result  $y_2^{(3)}$  due to nonzeros in  $\{a_{2,3}, a_{2,4}\}$  with  $\text{map}(a_{2,3}) = r_2$  and  $\text{map}(a_{2,4}) = c_4$ .

Similar to Section 3, after obtaining bipartitions  $\Pi(\mathcal{A}_k) = \{\mathcal{A}_{k,L}, \mathcal{A}_{k,R}\}$ ,  $\Pi(\mathcal{X}_k) = \{\mathcal{X}_{k,L}, \mathcal{X}_{k,R}\}$ , and  $\Pi(\mathcal{Y}_k) = \{\mathcal{Y}_{k,L}, \mathcal{Y}_{k,R}\}$ , the labels of the parts represented by the leaf nodes are updated in such a way that the resulting  $(K' + 1)$ -way partitions are denoted by  $\Pi_{K'+1}(\mathcal{A}) = \{\mathcal{A}_1, \dots, \mathcal{A}_{K'+1}\}$ ,  $\Pi_{K'+1}(\mathcal{X}) = \{\mathcal{X}_1, \dots, \mathcal{X}_{K'+1}\}$ , and  $\Pi_{K'+1}(\mathcal{Y}) = \{\mathcal{Y}_1, \dots, \mathcal{Y}_{K'+1}\}$ .

#### 4.1. Adaptation for conformal partitioning

Adapting the medium-grain model for a conformal partition on vectors  $x$  and  $y$  slightly differs from adapting the fine-grain model. Vertex set  $\mathcal{V}_k$  contains an amalgamated vertex  $v_i^{x/y}$  if at least one of the following conditions holds:

- $x_i \in \mathcal{X}_k$ , or equivalently,  $y_i \in \mathcal{Y}_k$ .
- $\exists a_{t,i} \in \mathcal{A}_k$  s.t.  $\text{map}(a_{t,i}) = c_i$ .
- $\exists a_{i,t} \in \mathcal{A}_k$  s.t.  $\text{map}(a_{i,t}) = r_i$ .

The weight of  $v_i$  is assigned as

$$w(v_i) = |\{a_{t,i} : a_{t,i} \in \mathcal{A}_k \text{ and } \text{map}(a_{t,i}) = c_i\}| + |\{a_{i,t} : a_{i,t} \in \mathcal{A}_k \text{ and } \text{map}(a_{i,t}) = r_i\}|.$$

Each volume/message net that connects  $v_i^x$  or  $v_i^y$  in  $\mathcal{H}_k^M$  now connects the amalgamated vertex  $v_i^{x/y}$ .

#### 5. Delayed addition and thresholding for message nets

Utilization of the message nets decreases the importance attributed to the volume nets in the partitioning process and this may lead to a relatively high bandwidth cost compared to the

case where no message nets are utilized. The more the number of RB steps in which the message nets are utilized, the higher the total communication volume. A high bandwidth cost can especially be attributed to the bipartitionings in the early levels of the RB tree. There are only a few nodes in the early levels of the RB tree compared to the late levels and each of these nodes represents a large processor group. The messages among these large processor groups are difficult to refrain from. In terms of hypergraph partitioning, since the message nets in the hypergraphs at the early levels of the RB tree connect more vertices and the cost of the message nets is much higher than the cost of the volume nets ( $t_{su} \gg t_w$ ), it is very unlikely for these message nets to be uncut. While the partitioner tries to save these nets from the cut in the early bipartitionings, it may cause high number of volume nets to be cut, which in turn are likely to introduce new messages in the late levels of the RB tree. Therefore, adding message nets in the early levels of the RB tree adversely affects the overall partition quality in multiple ways.

The RB approach provides the ability to adjust the partitioning parameters in the individual RB steps for the sake of the overall partition quality. In our model, we use this flexibility to exploit the trade-off between the bandwidth and latency costs by selectively deciding whether to add message nets in each bipartitioning. To make this decision, we use the level information of the RB steps in the RB tree. For a given  $L < \log K$ , the addition of the message nets is delayed until the  $L$ th level of the RB tree, i.e., the bipartitionings in level  $\ell$  are performed only with the volume nets for  $0 \leq \ell < L$ . Thus, the message nets are included in the bipartitionings in which they are expected to connect relatively fewer vertices.

Using a delay parameter  $L$  aims to avoid large message nets by not utilizing them in the early levels of the RB tree. However, there may still exist such nets in the late levels depending on the structure of the matrix being partitioned. Another idea is to eliminate the message nets whose size is larger than a given threshold. That is, for a given threshold  $T > 0$ , a message net  $n$  with  $|\text{Pins}(n)| > T$  is excluded from the corresponding bipartition. This approach also enables a selective approach for send and receive message nets. In our implementation of the row-column-parallel SpMV, the receive operations are performed by non-blocking MPI functions (i.e., `MPI_Irecv`), whereas the send operations are performed by blocking MPI functions (i.e., `MPI_Send`). When the maximum message count or the maximum communication volume is considered to be a serious bottleneck, blocking send operations may be more limiting compared to non-blocking receive operations. Note that saving message nets from the cut tends to assign the respective communication operations to fewer processors, hence the maximum message count and maximum communication volume may increase. Hence, a smaller threshold is preferable for the send message nets while a higher threshold is preferable for the receive nets.

#### 6. Experiments

We consider a total of five partitioning models for evaluation. Four of them are nonzero-based partitioning models: the fine-grain model (FG), the medium-grain model (MG), and the proposed models which simultaneously reduce the bandwidth and latency costs, as described in Section 3 (FG-LM) and Section 4 (MG-LM). The last partitioning model tested is the one-dimensional model (1D-LM) that simultaneously reduces the bandwidth and latency costs [21]. Two of these five models (FG and MG) encapsulate a single communication cost metric, i.e., total volume, while three of them (FG-LM, MG-LM, and 1D-LM) encapsulate two communication cost metrics, i.e., total volume and total message count. The partitioning constraint of balancing part weights in all these models corresponds to balancing of the computational loads of processors.



**Table 2**

The communication cost metrics obtained by the nonzero-based partitioning models with varying delay values ( $L$ ).

Model	$L$	Volume		Message	
		Max	Total	Max	Total
FG	–	567	52,357	60	5560
FG-LM	1	2700	96,802	56	2120
FG-LM	4	2213	94,983	49	2186
FG-LM	5	1818	90,802	46	2317
FG-LM	6	1346	82,651	46	2694
FG-LM	7	926	69,572	49	3574
MG	–	558	49,867	57	5103
MG-LM	1	1368	77,479	50	2674
MG-LM	4	1264	77,227	48	2735
MG-LM	5	1148	74,341	47	2809
MG-LM	6	969	69,159	47	3066
MG-LM	7	776	61,070	50	3695

In the models that address latency cost with the message nets, the cost of the volume nets is set to 1 while the cost of the message nets is set to 50, i.e., it is assumed  $t_{su} = 50t_w$ , which is also the setting recommended in [21].

The performance of the compared models are evaluated in terms of the partitioning cost metrics and the parallel SpMV run-time. The partitioning cost metrics include total volume, total message count, and load imbalance (these are explained in detail in following sections) and they are helpful to test the validity of the proposed models. At each RB step in all models, we used PaToH [7] in the default settings to obtain a two-way partition of the respective hypergraph. An imbalance ratio of 10% is used in all models, i.e.,  $\epsilon = 0.10$ . We test for five different number of parts/processors,  $K \in \{64, 128, 256, 512, 1024\}$ . The parallel SpMV is implemented using the PETSc toolkit [3]. PETSc contains structures and routines for parallel solution of applications modeled by partial differential equations. It supports MPI-based and hybrid parallelism, and offers a wide range of sparse linear solvers and preconditioners. The parallel SpMV realized within PETSc is run on a Blue Gene/Q system using the partitions provided by the five compared models. A node on Blue Gene/Q system consists of 16 PowerPC A2 processors with 1.6 GHz clock frequency and 16 GB memory.

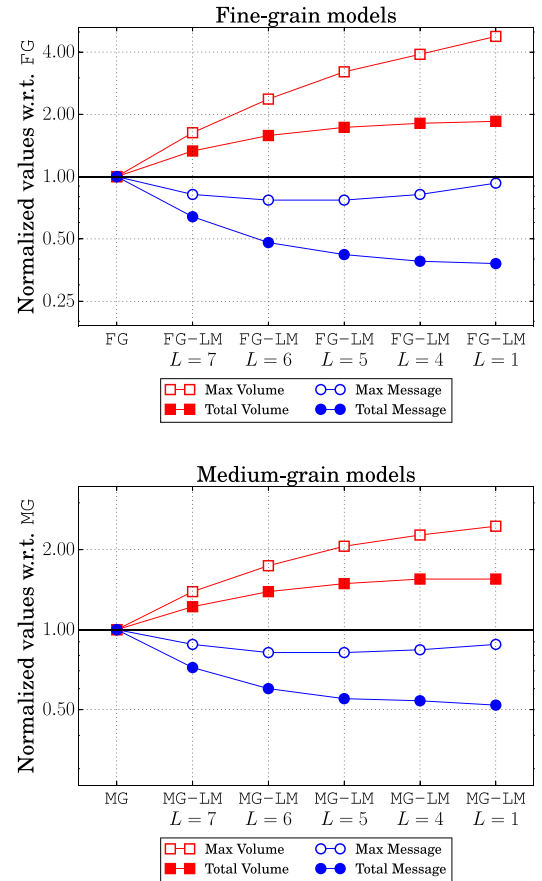
The experiments are performed on an extensive dataset containing matrices from the SuiteSparse Matrix Collection [11]. We consider the case of conformal vector partitioning as it is more common for the applications in which SpMV is used as a kernel operation. Hence, only the square matrices are considered. We use the following criteria for the selection of test matrices: (i) the minimum and maximum number of nonzeros per processor are, respectively, set to 100 and 100,000, (ii) the matrices that have more than 50 million nonzeros are excluded, and (iii) the minimum number of rows/columns per processor is set to 50. The resulting number of matrices are 833, 730, 616, 475, and 316 for  $K = 64, 128, 256, 512$ , and 1024 processors, respectively. The union of these sets of matrices makes up to a total of 978 matrices.

### 6.1. Tuning parameters for nonzero-based partitioning models

There are two important issues described in Section 5 regarding the addition of the message nets for the nonzero-based partitioning models. We next discuss setting these parameters.

#### 6.1.1. Delay parameter ( $L$ )

We investigate the effect of the delay parameter  $L$  on four different communication cost metrics for the fine-grain and medium-grain models with the message nets. These cost metrics are maximum volume, total volume, maximum message count, and total message count. The volume metrics are in terms of number of words communicated. We compare FG-LM with delay against FG,



**Fig. 7.** The effect of the delay parameter on nonzero-based partitioning models in four different communication metrics.

as well as MG-LM with delay against MG. We only present the results for  $K = 256$  since the observations made for the results of different  $K$  values are similar. Note that there are  $\log 256 = 8$  bipartitioning levels in the corresponding RB tree. The tested values of the delay parameter  $L$  are 1, 4, 5, 6, and 7. Note that the message nets are added in a total of 4, 3, 2, and 1 levels for the  $L$  values of 4, 5, 6, and 7, respectively. When  $L = 1$ , it is equivalent to adding message nets throughout the whole partitioning without any delay. Note that it is not possible to add message nets at the root level (i.e., by setting  $L = 0$ ) since there is no partition available yet to form the message nets. The results for the remaining values of  $L$  are not presented, as the tested values contain all the necessary insight for picking a value for  $L$ . Table 2 presents the results obtained. The value obtained by a partitioning model for a specific cost metric is the geometric mean of the values obtained for the matrices by that partitioning model (i.e., the mean of the results for 616 matrices). We also present two plots in Fig. 7 to provide a visual comparison of the values presented in Table 2. The plot at the top belongs to the fine-grain models and each different cost metric is represented by a separate line in which the values are normalized with respect to those of the standard fine-grain model FG. Hence, a point on a line below  $y = 1$  indicates the variants of FG-LM attaining a better performance in the respective metric compared to FG, whereas a point in a line above indicates a worse performance. For example, FG-LM with  $L = 7$  attains 0.72 times the total message count of FG, which corresponds to the second point of the line marked with a filled circle. The plot at the bottom compares the medium-grain models in a similar fashion.

It can be seen from Fig. 7 that, compared to FG, FG-LM attains better performance in maximum and total message count, and a

**Table 3**

The communication cost metrics of FG-LM with varying message net thresholds ( $T_S$ ,  $T_R$ ).

$T_S$	$T_R$	Volume		Message	
		Max	Total	Max	Total
–	–	1346	82,651	46	2694
15	15	706	56,218	58	4539
15	30	773	58,452	56	4258
15	50	835	60,864	54	4043
30	15	793	58,418	59	4251
30	30	827	60,086	57	4087
30	50	900	62,393	55	3879
50	15	879	61,099	59	4037
50	30	908	62,516	58	3877
50	50	952	64,041	56	3729

worse performance in maximum and total volume. A similar observation is also valid for comparing MG with MG-LM. As the number of RB tree levels in which the message nets are added increases, FG-LM and MG-LM obtain lower latency and higher bandwidth overheads compared to FG and MG, respectively. The improvement rates in latency cost obtained by the partitioning models utilizing the message nets saturate around  $L = 6$  or  $L = 5$ , whereas the deterioration rates in bandwidth cost continue to increase. In other words, adding message nets in the bipartitionings other than those in the last two or three levels of the RB tree has small benefits in terms of improving the latency cost but it has a substantial negative effect on the bandwidth cost, especially on maximum volume. For this reason, we choose FG-LM and MG-LM with  $L = 6$ , i.e., add message nets in the last two levels of the RB tree.

#### 6.1.2. Message net threshold parameters ( $T_S$ , $T_R$ )

The message net threshold parameters for the send and receive message nets are, respectively, denoted with  $T_S$  and  $T_R$ . The tested values are set based upon the average degree of the message nets throughout the partitioning, which is found to be close to 30. We evaluate threshold values smaller than, roughly equal to, and greater than this average degree:  $T_S, T_R \in \{15, 30, 50\}$ . We follow a similar experimental setting as for the delay parameter and only present the results for  $K = 256$ . In addition, we omit the discussions for the medium-grain models as the observations made for the fine-grain and medium-grain models are alike. Table 3 presents the values for four different cost metrics obtained by FG-LM and

FG-LM with nine different threshold settings. Note that the delay value of  $L = 6$  is utilized in all these experiments.

The partitionings without large message nets lead to lower bandwidth and higher latency costs as seen in Table 3 compared to the case without any threshold, i.e., FG-LM. The more the number of eliminated message nets, the higher the latency cost and the lower the bandwidth cost. Among the nine combinations for  $T_S$  and  $T_R$  in the table, we pick  $T_S = 15$  and  $T_R = 50$  due to its reasonable maximum volume and maximum message count values for the reasons described in Section 5.

### 6.2. Comparison of all partitioning models

#### 6.2.1. Partitioning cost metrics

We present the values obtained by the four nonzero-based partitioning models in six different partitioning cost metrics in Table 4. These cost metrics are computational imbalance (indicated in the column titled “Imb (%)”), maximum and total volume, maximum and total message count, and partitioning time in seconds. Note that these metrics include the three metrics described in Section 2.2. Each entry in the table is the geometric mean of the values for the matrices that belong to the respective value of  $K$ . The columns three to eight in the table display the actual values, whereas the columns nine to fourteen display the normalized values, where the results obtained by FG-LM and MG-LM at each  $K$  value are normalized with respect to those obtained by FG and MG at that  $K$  value, respectively. The top half of the table displays the results obtained by the fine-grain models, whereas the bottom half displays the results obtained by the medium-grain models.

Among the four nonzero-based partitioning models compared in Table 4, the models that consider both the bandwidth and latency overheads achieve better total and maximum message counts compared to the models that solely consider the bandwidth overhead. For example at  $K = 256$ , FG-LM attains 27% improvement in total message count compared to FG, while MG-LM attains 24% improvement in total message count compared to MG. On the other hand, the two models that solely consider the bandwidth overhead achieve better total and maximum volume compared to the two models that also consider the latency overhead. This is because FG and MG optimize a single cost metric, while FG-LM and MG-LM aim to optimize two cost metrics at once. At  $K = 256$ , FG-LM causes 16% deterioration in total volume compared to FG, while MG-LM causes 18% deterioration in total volume compared to

**Table 4**

Comparison of nonzero-based partitioning models in six cost metrics.

$K$	Model	Actual values						Normalized values w.r.t. FG/MG					
		Imb (%)	Volume		Message		Part. time	Imb	Volume		Message		Part. time
			Max	Total	Max	Total			Max	Total	Max	Total	
64	FG	0.91	413	11,811	32	968	7.7	–	–	–	–	–	–
	FG-LM	0.88	542	13,267	29	753	7.4	0.97	1.31	1.12	0.91	0.78	0.97
128	FG	1.11	484	24,670	45	2,332	16.4	–	–	–	–	–	–
	FG-LM	1.01	669	28,159	40	1,751	16.3	0.91	1.38	1.14	0.89	0.75	1.00
256	FG	1.36	567	52,357	60	5,560	40.9	–	–	–	–	–	–
	FG-LM	1.21	835	60,864	54	4,043	40.8	0.89	1.47	1.16	0.90	0.73	1.00
512	FG	1.67	584	92,141	72	11,186	77.9	–	–	–	–	–	–
	FG-LM	1.61	863	108,497	66	8,218	77.2	0.96	1.48	1.18	0.92	0.73	0.99
1024	FG	1.87	530	165,923	69	20,209	156.2	–	–	–	–	–	–
	FG-LM	1.81	811	196,236	66	15,415	159.6	0.97	1.53	1.18	0.96	0.76	1.02
64	MG	0.90	412	11,655	31	928	3.9	–	–	–	–	–	–
	MG-LM	0.87	521	13,205	28	732	4.1	0.97	1.26	1.13	0.90	0.79	1.06
128	MG	1.13	482	24,256	44	2,217	8.1	–	–	–	–	–	–
	MG-LM	1.08	634	27,799	39	1,690	8.4	0.96	1.32	1.15	0.89	0.76	1.04
256	MG	1.48	558	49,867	57	5,103	19.1	–	–	–	–	–	–
	MG-LM	1.39	766	58,981	52	3,876	20.6	0.94	1.37	1.18	0.91	0.76	1.08
512	MG	1.91	588	91,856	67	10,265	39.7	–	–	–	–	–	–
	MG-LM	1.80	785	108,128	62	7,878	43.7	0.94	1.34	1.18	0.93	0.77	1.10
1024	MG	2.05	530	165,722	65	18,692	82.2	–	–	–	–	–	–
	MG-LM	2.00	724	196,443	61	14,827	87.5	0.98	1.37	1.19	0.94	0.79	1.06

**Table 5**Comparison of partitioning models in six cost metrics at  $K = 256$ .

Model	Imb (%)	Volume		Message		Part. time
		Max	Total	Max	Total	
1D-LM	2.50	968	101,565	33	2448	13.2
FG	1.36	567	52,357	60	5560	40.9
FG-LM	1.21	835	60,864	54	4043	40.8
MG	1.48	558	49,867	57	5103	19.1
MG-LM	1.39	766	58,981	52	3876	20.6

MG. Note that the models behave accordingly in maximum volume and maximum message count metrics as although these metrics are not directly addressed by any of the models, the former one is largely dependent on the total volume while the latter one is largely dependent on the total message count. FG-LM and MG-LM have slightly lower imbalance compared to FG and MG, respectively. Addition of the message nets does not seem to change the partitioning overhead, a result likely to be a consequence of the choice of the delay and net threshold parameters.

Another observation worth discussion is the performance of the medium-grain models against the performance of the fine-grain models. When MG is compared to FG or MG-LM is compared to FG-LM, the medium-grain models achieve slightly better results in volume and message cost metrics, and slightly worse results in imbalance. However, the partitioning overhead of the medium-grain models is much lower than the partitioning overhead of the fine-grain models: the medium grain models are 1.8–2.2 $\times$  faster. This is also one of the main findings of [19], which makes the medium-grain model a better alternative for obtaining nonzero-based partitions.

1D-LM and nonzero-based partitioning models are compared in Table 5 at  $K = 256$ . 1D-LM has higher total volume and imbalance, and lower total message count compared to the nonzero-based partitioning models. The nonzero-based models have broader search space due to their representation of the SpMV via smaller units, which allows them to attain better volume and imbalance. The latency overheads of FG and MG are higher than the latency overhead of 1D-LM simply because latency is not addressed in the former two. Although FG-LM and MG-LM may as well obtain comparable latency overheads with 1D-LM (e.g., compare total message count of FG-LM with  $L = 1$  in Table 2 against total message count of 1D-LM in Table 5), we favor a decrease in volume-related cost metrics at the expense of a small deterioration in latency-related cost metrics in these two models. 1D-LM has the lowest partitioning overhead due to having the smallest hypergraph among the five models. A similar discussion follows for the maximum volume and maximum message count metrics as for the total volume and total message count metrics.

In the remainder of the paper, we use MG and MG-LM among the nonzero-based models for evaluation due to their lower partitioning overhead and slightly better performance compared to FG and FG-LM, respectively, in the remaining metrics.

### 6.2.2. Parallel SpMV performance

We compare 1D-LM, MG, and MG-LM in terms of parallel SpMV runtime. Parallel SpMV is run with the partitions obtained through these three models. There are 12 matrices tested, listed with their types as follows: eu-2005 (web graph), ford2 (mesh), Freescale1 (circuit simulation), invextr1\_new (computational fluid dynamics), k1\_san (2D/3D), LeGresley\_87936 (power network), mouse\_gene (gene network), olesnik0 (2D/3D), tuma1 (2D/3D), turon\_m (2D/3D), usroads (road network), web-Google (web graph). Number of nonzeros in these matrices varies between 87,760 and 28,967,291. These 12 matrices constitute a subset of 978 matrices for which the partitioning models are

compared in terms of partitioning cost metrics in the preceding sections. Four different number of processors (i.e.,  $K$ ) are tested: 64, 128, 256, and 512. We did not test for 1024 processors as in most of the tested matrices SpMV could not scale beyond 512 processors. We only consider the strong-scaling case. The parallel SpMV is run for 100 times and the average runtime (in milliseconds) is reported. The obtained results are presented in Fig. 8.

The plots in Fig. 8 show that both MG and MG-LM scale usually better than 1D-LM. It is known the nonzero-based partitioning models scale better than the 1D models due to their lower communication overheads and computational imbalance. In difficult instances such as invextr1\_new or mouse\_gene at which 1D-LM does not scale, using a nonzero-based model such as MG or MG-LM successfully scales the parallel SpMV. MG-LM improves the scalability of MG in most of the test instances. Apart from the instances Freescale1, invextr1\_new, and turon\_m, MG-LM performs significantly better than MG. MG-LM's performance especially gets more prominent with increasing number of processors, which is due to the fact that the latency overheads are more critical in the overall communication costs in high processor counts since the message size usually decreases with increasing number of processors. These plots show that using a nonzero-based partitioning model coupled with the addressing of multiple communication cost metrics yields the best parallel SpMV performance.

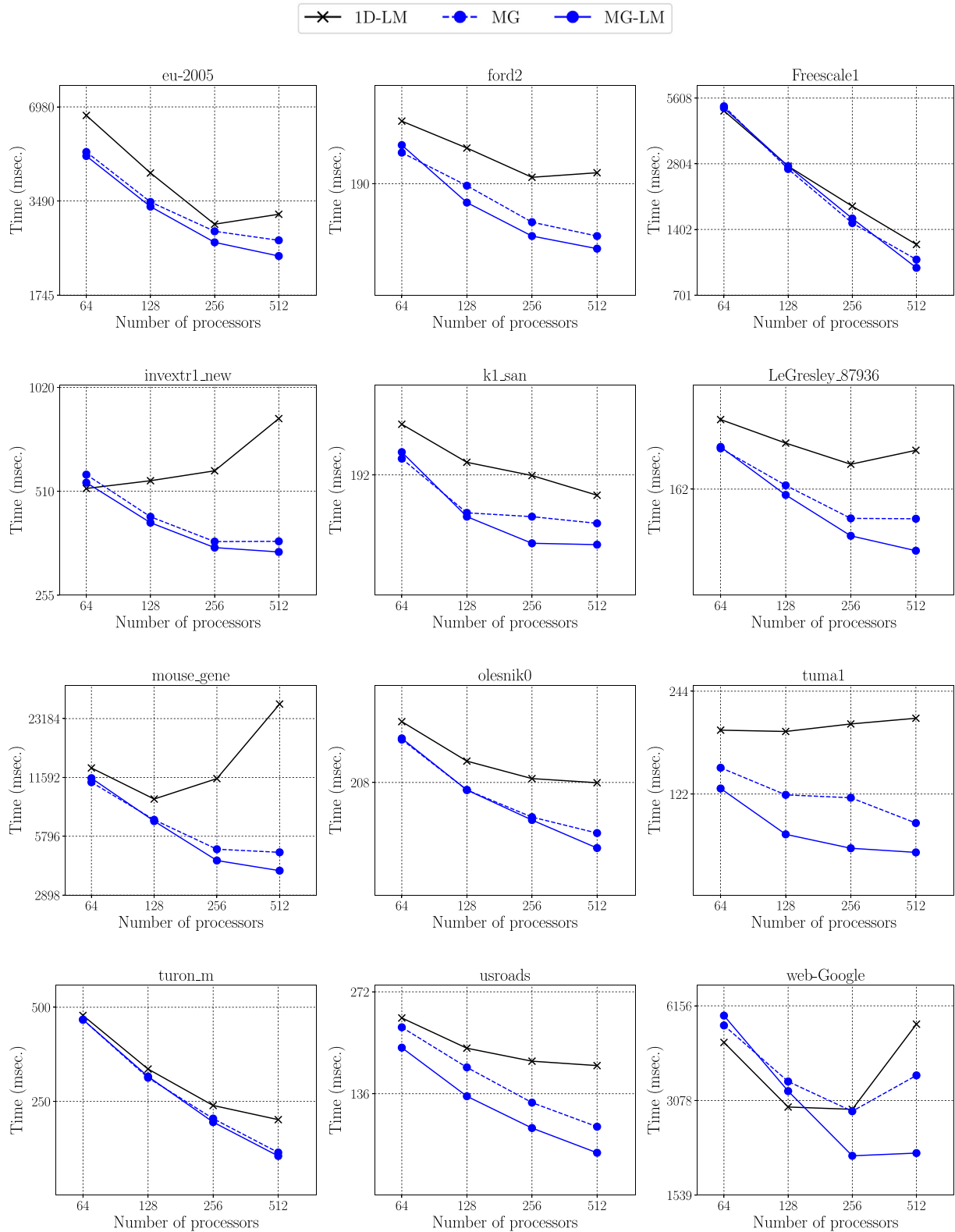
## 7. Related work

We first provide a brief taxonomy of the sparse matrix partitioning models for parallel SpMV in order to place the proposed models in their context by arguing their pros and cons. Then, we describe the works that have similar goals with this work, i.e., the reduction of multiple cost metrics, of which at least one is a latency-related cost metric.

The literature is rich in terms of the sparse matrix partitioning models based on graph partitioning [5,14,16,20] and hypergraph partitioning (HP) [1,2,7–10,12,17–19,21–25]. The HP models are grouped into two as one-dimensional (1D) and two-dimensional (2D) depending on the dimensions of the matrix partitioning. Fig. 9 provides a visual taxonomy for the partitioning models mentioned in this section. 1D models [7] partition the nonzeros in the matrix along a single dimension, either row or column, while 2D models partition them along both row and column dimensions. The two main variants of the 1D partitioning are the row-wise and column-wise partitionings. The 2D models are generally superior to the 1D models in terms of parallel SpMV performance due to their higher flexibility in distributing the matrix nonzeros. The three main variants of the 2D partitioning are the checkerboard [9,10], jagged [10], and nonzero-based [8,10,19] partitionings. A 2D virtual processor mesh is assumed for the checkerboard and jagged partitionings in order to provide upper bounds on communication overheads.

Unlike the partitioning models mentioned so far, the nonzero-based partitioning does not possess any structure in terms of the nonzero distribution among processors, so, it is the most general partitioning possible [24]. The two variants of the nonzero-based partitioning are the fine-grain and medium-grain partitionings, in which each nonzero and groups of nonzeros are treated individually during the partitioning, respectively. This work focuses on improving the fine-grain and medium-grain HP models and they differ from the HP models proposed in [8,19] for obtaining nonzero-based partitionings in the sense that our models address both the bandwidth and latency costs, while they address only the bandwidth cost.

The existing partitioning models that address the bandwidth and latency costs in the literature can be grouped into two according to whether they explicitly address the latency cost or not (the bandwidth cost is usually addressed explicitly). The models



**Fig. 8.** Comparison of partitioning models in terms of parallel SpMV runtime.

that do not explicitly address the latency cost provide an upper bound on the message counts [5,9,10]. We focus on the works that explicitly address the latency cost [12,21,22], which is also the case in this work. Among these works, the one proposed in [22] is a two-phase approach which addresses the bandwidth cost in the first phase utilizing a 1D model and the latency cost in the second phase

utilizing the communication hypergraph model. In the two-phase approaches, since different cost metrics are addressed in separate phases, a metric minimized in a particular phase may deteriorate in the other phase. Our models fall into the category of single-phase approaches. The other two works also adopt a single-phase approach to address multiple communication cost metrics, where



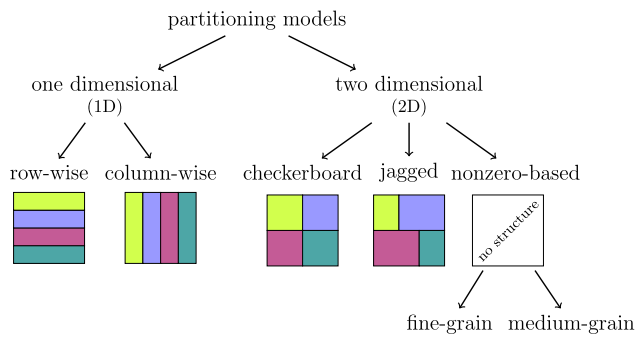


Fig. 9. Partitioning taxonomy.

UMPa [12] uses a direct  $K$ -way partitioning approach, while [21] exploits the recursive bipartitioning paradigm. UMPa is rather expensive as it introduces an additional cost involving a quadratic factor in terms of the number of processors to each refinement pass. Our approach introduces an additional cost involving a mere logarithmic factor in terms of the number of processors to the entire partitioning. Furthermore, our approach enables the use of existing HP tools for bipartitioning within the RB framework.

The message net concept has recently been proposed for certain types of iterative applications that involve a computational phase either preceded or followed by a communication phase with a restriction of conformal partitions on input and output data [21]. 1D row-parallel and column-parallel SpMV operations constitute examples for these applications. This work differs from [21] in the sense that the nonzero-based partitions necessitate a parallel SpMV that involves two communication phases with no restriction of conformal partitions on input and output vectors of the SpMV operation. Partitions on the vectors can also be found after finding a partition on matrix nonzeros [4,22]. Our proposed HP models, on the other hand, find all three partitions at once in a single partitioning phase.

## 8. Conclusion

We proposed two novel nonzero-based matrix partitioning models, a fine-grain and a medium-grain model, that simultaneously address the bandwidth and latency costs of parallel SpMV. These models encapsulate two communication cost metrics at once as opposed to their existing counterparts which only address a single cost metric regarding the bandwidth cost. Our approach exploits the recursive bipartitioning paradigm to incorporate the latency minimization into the partitioning objective via message nets. In addition, we proposed two practical enhancements to find a good balance between reducing the bandwidth and the latency costs. The experimental results obtained on an extensive dataset show that the proposed models attain up to 27% improvement in latency-related cost metrics over their existing counterparts on average and the scalability of parallel SpMV can substantially be improved with the proposed models.

## Acknowledgment

We acknowledge PRACE for awarding us access to resources Juqueen (Blue Gene/Q) based in Germany at Jülich Supercomputing Centre.

## References

- [1] S. Acer, O. Selvitopi, C. Aykanat, Improving performance of sparse matrix dense matrix multiplication on large-scale parallel systems, in: *Theory and Practice of Irregular Applications*, Parallel Comput. 59 (2016) 71–96.
- [2] K. Akbudak, E. Kayaaslan, C. Aykanat, Hypergraph partitioning based models and methods for exploiting cache locality in sparse matrix-vector multiplication, *SIAM J. Sci. Comput.* 35 (3) (2013).
- [3] S. Balay, S. Abhyankar, M.F. Adams, J. Brown, P. Brune, K. Buschelman, V. Eijkhout, W.D. Gropp, D. Kaushik, M.G. Knepley, L.C. McInnes, K. Rupp, B.F. Smith, H. Zhang, PETSc Users Manual, Tech. Rep. ANL-95/11 - Revision 3.5, Argonne National Laboratory, 2014, URL <http://www.mcs.anl.gov/petsc>.
- [4] R.H. Bisseling, W. Meesen, Communication balancing in parallel sparse matrix-vector multiply, *Electron. Trans. Numer. Anal.* 21 (2005) 47–65.
- [5] E.G. Boman, K.D. Devine, S. Rajamanickam, Scalable matrix computations on large scale-free graphs using 2D graph partitioning, in: *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*, SC '13, ACM, New York, NY, USA, 2013, pp. 50:1–50:12, <http://dx.doi.org/10.1145/2503210.2503293>, URL <http://doi.acm.org/10.1145/2503210.2503293>.
- [6] U. Çatalyürek, C. Aykanat, Decomposing irregularly sparse matrices for parallel matrix-vector multiplication, in: *Proceedings of the Third International Workshop on Parallel Algorithms for Irregularly Structured Problems*, IR-REGULAR '96, Springer-Verlag, London, UK, 1996, pp. 75–86, URL <http://dl.acm.org/citation.cfm?id=646010.676990>.
- [7] U.V. Çatalyürek, C. Aykanat, Hypergraph-partitioning-based decomposition for parallel sparse-matrix vector multiplication, *IEEE Trans. Parallel Distrib. Syst.* 10 (7) (1999) 673–693, <http://dx.doi.org/10.1109/71.780863>.
- [8] U. Çatalyürek, C. Aykanat, A fine-grain hypergraph model for 2D decomposition of sparse matrices, in: *Proceedings of the 15th International Parallel and Distributed Processing Symposium*, IPDPS '01, IEEE Computer Society, Washington, DC, USA, 2001, p. 118, URL <http://dl.acm.org/citation.cfm?id=645609.663255>.
- [9] U. Çatalyürek, C. Aykanat, A hypergraph-partitioning approach for coarse-grain decomposition, in: *Proceedings of the 2001 ACM/IEEE Conference on Supercomputing*, SC '01, ACM, New York, NY, USA, 2001, <http://dx.doi.org/10.1145/582034.582062>, pp. 28–28, URL <http://doi.acm.org/10.1145/582034.582062>.
- [10] U.V. Çatalyürek, C. Aykanat, B. Uçar, On two-dimensional sparse matrix partitioning: models, methods, and a recipe, *SIAM J. Sci. Comput.* 32 (2) (2010) 656–683, <http://dx.doi.org/10.1137/080737770>.
- [11] T.A. Davis, Y. Hu, The University of Florida sparse matrix collection, *ACM Trans. Math. Software* 38 (1) (2011) 1:1–1:25, <http://dx.doi.org/10.1145/2049662.2049663>, URL <http://doi.acm.org/10.1145/2049662.2049663>.
- [12] M. Deveci, K. Kaya, B. Uçar, Ü. Çatalyürek, Hypergraph partitioning for multiple communication cost metrics: Model and methods, *J. Parallel Distrib. Comput.* 77 (2015) 69–83, <http://dx.doi.org/10.1016/j.jpdc.2014.12.002>, URL <http://www.sciencedirect.com/science/article/pii/S0743731514002275>.
- [13] B. Hendrickson, Graph partitioning and parallel solvers: Has the emperor no clothes? (extended abstract), in: *Proceedings of the 5th International Symposium on Solving Irregularly Structured Problems in Parallel*, IRREGULAR '98, Springer-Verlag, London, UK, 1998, pp. 218–225, URL <http://dl.acm.org/citation.cfm?id=646012.677019>.
- [14] B. Hendrickson, T.G. Kolda, Partitioning rectangular and structurally unsymmetric sparse matrices for parallel processing, *SIAM J. Sci. Comput.* 21 (6) (1999) 2048–2072, <http://dx.doi.org/10.1137/S1064827598341475>, URL <http://dx.doi.org/10.1137/S1064827598341475>.
- [15] B. Hendrickson, T.G. Kolda, Graph partitioning models for parallel computing, *Parallel Comput.* 26 (12) (2000) 1519–1534, [http://dx.doi.org/10.1016/S0167-8191\(00\)00048-X](http://dx.doi.org/10.1016/S0167-8191(00)00048-X).
- [16] G. Karypis, V. Kumar, A fast and high quality multilevel scheme for partitioning irregular graphs, *SIAM J. Sci. Comput.* 20 (1) (1998) 359–392, <http://dx.doi.org/10.1137/S1064827595287997>.
- [17] E. Kayaaslan, C. Aykanat, B. Uçar, 1.5D parallel sparse matrix-vector multiply, *SIAM J. Sci. Comput.* 40 (1) (2018) C25–C46, <http://dx.doi.org/10.1137/16M1105591>.
- [18] E. Kayaaslan, B. Uçar, C. Aykanat, Semi-two-dimensional partitioning for parallel sparse matrix-vector multiplication, in: *Parallel and Distributed Processing Symposium Workshop (IPDPSW)*, 2015 IEEE International, 2015, pp. 1125–1134, <http://dx.doi.org/10.1109/IPDPSW.2015.20>.
- [19] D. Pelt, R. Bisseling, A medium-grain method for fast 2D bipartitioning of sparse matrices, in: *Parallel and Distributed Processing Symposium*, 2014 IEEE 28th International, 2014, 529–539, <http://dx.doi.org/10.1109/IPDPS.2014.62>.
- [20] K. Schloegel, G. Karypis, V. Kumar, Parallel multilevel algorithms for multi-constraint graph partitioning, in: A. Bode, T. Ludwig, W. Karl, R. Wismüller (Eds.), *Euro-Par 2000 Parallel Processing*, in: *Lecture Notes in Computer Science*, vol. 1900, Springer Berlin Heidelberg, 2000, pp. 296–310, [http://dx.doi.org/10.1007/3-540-44520-X\\_39](http://dx.doi.org/10.1007/3-540-44520-X_39), URL [http://dx.doi.org/10.1007/3-540-44520-X\\_39](http://dx.doi.org/10.1007/3-540-44520-X_39).
- [21] O. Selvitopi, S. Acer, C. Aykanat, A recursive hypergraph bipartitioning framework for reducing bandwidth and latency costs simultaneously, *IEEE Trans. Parallel Distrib. Syst.* 28 (2) (2017) 345–358, <http://dx.doi.org/10.1109/TPDS.2016.2577024>.

- [22] B. Uçar, C. Aykanat, Encapsulating multiple communication-cost metrics in partitioning sparse rectangular matrices for parallel matrix-vector multiplies, *SIAM J. Sci. Comput.* 25 (6) (2004) 1837–1859, <http://dx.doi.org/10.1137/S1064827502410463>.
- [23] B. Uçar, C. Aykanat, Revisiting hypergraph models for sparse matrix partitioning, *SIAM Rev.* 49 (2007) 595–603, <http://dx.doi.org/10.1137/060662459>, URL <http://portal.acm.org/citation.cfm?id=1330215.1330219>.
- [24] B. Vastenhouw, R.H. Bisseling, A two-dimensional data distribution method for parallel sparse matrix-vector multiplication, *SIAM Rev.* 47 (2005) 67–95, <http://dx.doi.org/10.1137/S0036144502409019>, URL <http://portal.acm.org/citation.cfm?id=1055334.1055397>.
- [25] A.N. Yzelman, R.H. Bisseling, Cache-oblivious sparse matrix-vector multiplication by using sparse matrix partitioning methods, *SIAM J. Sci. Comput.* 31 (4) (2009) 3128–3154, <http://dx.doi.org/10.1137/080733243>.



**Seher Acer** received her B.S., M.S. and Ph.D. degrees in Computer Engineering from Bilkent University, Turkey, where she is currently a postdoctoral researcher. Her research interests are combinatorial scientific computing, graph and hypergraph partitioning for sparse matrix and tensor computations, and parallel computing.



**Oguz Selvitopi** received his B.S., M.S. and Ph.D. degrees in Computer Engineering from Bilkent University, Turkey, where he is currently a postdoctoral researcher. His research interests are parallel computing, scientific computing, and parallel and distributed systems.



**Cevdet Aykanat** received the BS and MS degrees from Middle East Technical University, Ankara, Turkey, both in Electrical Engineering, and the Ph.D. degree from Ohio State University, Columbus, in Electrical and Computer Engineering. He worked at the Intel Supercomputer Systems Division, Beaverton, Oregon, as a research associate. Since 1989, he has been affiliated with the Department of Computer Engineering, Bilkent University, Ankara, Turkey, where he is currently a professor. His research interests mainly include parallel computing, parallel scientific computing and its combinatorial aspects. He is the recipient of the 1995 Young Investigator Award of The Scientific and Technological Research Council of Turkey and 2007 Parlar Science Award. He has served as an Associate Editor of IEEE Transactions of Parallel and Distributed Systems between 2008 and 2012.