

Lawrence Berkeley National Laboratory

LBL Publications

Title

Simultaneous Computational and Data Load Balancing in Distributed-Memory Setting

Permalink

<https://escholarship.org/uc/item/3q45g481>

Journal

SIAM Journal on Scientific Computing, 44(6)

ISSN

1064-8275

Authors

Çeliktuğ, Mestan Firat
Karsavuran, M Ozan
Acer, Seher
[et al.](#)

Publication Date

2022-12-01

DOI

10.1137/22m1485772

Copyright Information

This work is made available under the terms of a Creative Commons Attribution-NonCommercial License, available at <https://creativecommons.org/licenses/by-nc/4.0/>

Peer reviewed

SIMULTANEOUS COMPUTATIONAL AND DATA LOAD BALANCING IN DISTRIBUTED-MEMORY SETTING

MESTAN FIRAT ÇELIKTUĞ^{*}, M. OZAN KARSAVURAN[†], SEHER ACER[‡], AND
CEVDET AYKANAT[†]

Abstract. Several successful partitioning models and methods have been proposed and used for computational load balancing of irregularly sparse applications in distributed-memory setting. However, the literature lacks partitioning models and methods that encode both computational and data load balancing. In this article, we try to close this gap in the literature by proposing two hypergraph partitioning (HP) models which simultaneously encode computational and data load balancing. Both models utilize a two-constraint formulation where the first constraint encodes the computational load and the second constraint encodes the data load. In the first model, we introduce explicit data vertices for encoding data load and we replicate those data vertices at each recursive bipartitioning (RB) step for encoding data replication. In the second model, we introduce a data weight distribution scheme for encoding data load and we update those weights at each RB step. The nice property of both proposed models is that they do not necessitate developing a new partitioner from scratch. Both models can easily be implemented by invoking any HP tool that supports multi-constraint partitioning as a two-way partitioner at each RB step. The validity of the proposed models is tested on two widely-used irregularly sparse applications: parallel mesh simulations and parallel sparse matrix sparse matrix multiplication (SpGEMM). Both proposed models achieve significant improvement over a baseline model.

Key words. computational load balance, data load balance, distributed-memory systems, hypergraph partitioning, recursive bipartitioning, multi-constraint partitioning, general sparse matrix-matrix multiplication, mesh partitioning

AMS subject classifications. 05C85,05C65,65F50,68R10

1. Introduction. In a distributed-memory setting, task-to-processor assignment has a significant role to attain high performance. This assignment affects multiple performance metrics such as computational load balance and communication costs which include bandwidth and latency components. Over the years, these metrics are widely studied alone or in a combination [1–3, 9–11, 14, 17, 20, 24, 36, 38, 39, 41–43]. There also exist combinatorial models and works which target minimizing these metrics under a given partial or complete data partition [6, 13, 18, 19, 21, 33]. In the high performance computing community, whenever load balance is pronounced it is almost always computational load balance [11, 12, 14, 18, 25, 25–27, 37]. In the cloud computing community, data load balance is considered, however, in that context data load balance is usually the only objective of the partitioning [30, 32, 35]. That is, only data is partitioned across data centers without associated tasks. In the literature, data load is considered for data migration cost [12, 19, 25] and memory capacity [5, 22, 40]. To our knowledge, this is the first paper in which simultaneous balance on computational and data loads are considered.

The target problem consists of atomic tasks and data elements to be assigned to the processors. Tasks do not have any computational dependency, whereas a data element might be needed by multiple tasks. If such tasks are assigned to different processors, then that data element will be replicated to those processors. Tasks are

^{*}Department of Computer Science, University of Texas at Dallas, and Department of Computer Engineering, Gazi University, Ankara, Turkey (mestanfirat.celiktug@utdallas.edu, fi-rat.celiktug@gazi.edu.tr).

[†]Department of Computer Engineering, Bilkent University, 06800 Ankara, Turkey (ozan.karsavuran@cs.bilkent.edu.tr, aykanat@cs.bilkent.edu.tr).

[‡]National Center for Computational Sciences, Oak Ridge National Laboratory (acers@ornl.gov).

44 associated with computational weights, whereas data elements are associated with
 45 sizes. Then the problem is finding the task-to-processor assignment in which maxi-
 46 mum computational and data loads of processors are minimized simultaneously.

47 In this work, we propose two novel hypergraph partitioning (HP) based models
 48 which simultaneously consider computational and data loads of the processors in a
 49 distributed-memory setting. In the first model, there exist vertices representing com-
 50 putational tasks as well as data elements. There exist nets also representing data
 51 elements and their relations with tasks. A two-constraint formulation is utilized for
 52 simultaneous computational and data load balancing, where partitioning objective
 53 encodes minimization of the total data replication. We utilize the recursive biparti-
 54 tioning (RB) scheme to increase the performance of the proposed vertex replication
 55 scheme by applying it at each level. For the RB framework, we propose a novel data
 56 vertex replication scheme to encode better data load balancing in the further RB
 57 steps.

58 In the second model, vertices represent computational tasks, whereas nets rep-
 59 resent data elements as well as their relations with tasks. In this model, a similar
 60 two-constraint formulation is also utilized. In order to model data loads, in contrast
 61 to the first model which contains explicit data vertices, we propose a data weight
 62 distribution to the vertex weights. Furthermore, we also utilize the RB framework for
 63 enabling utilization of the proposed weight distribution at each level.

64 In two-constraint formulation utilized in both models, finding balance on the com-
 65 putational loads through part weights encapsulates minimizing the computational load
 66 of the maximally loaded processor, since the total computational load is fixed. This is
 67 not true for the data load, because the total data load depends on the computational
 68 task partitioning. On the other hand, both proposed models minimize the amount
 69 of the data replication with the partitioning objective. In that way, proposed two-
 70 constraint formulation also encapsulates minimizing the data load of the maximally
 71 loaded processor.

72 We evaluate the performance of the proposed models against a HP based baseline
 73 model which only tries to balance computational loads of the processors. We utilized
 74 two sample applications in which our target problem arises: Parallel Finite Element
 75 Method and Volume Element Method based simulations which involve partitioning
 76 irregular 2D or 3D meshes and row-row-parallel Sparse Generalized Matrix Matrix
 77 Multiplication which involves partitioning two irregularly sparse input matrices. Ex-
 78 tensive experiments conducted for a wide range of partitioning instances on up to 1024
 79 processors show that both proposed models achieve significantly better performance
 80 than the baseline model.

81 The rest of the paper is organized as follows: [Section 2](#) describes the framework
 82 and formally defines the target problem. We give preliminary information about HP
 83 and RB framework in [section 3](#). We propose two HP models for the target problem in
 84 [section 4](#). In [section 5](#) experimental results are presented and discussed. We briefly
 85 discuss related works in [section 6](#). Finally, [section 7](#) concludes the paper.

86 **2. Framework and problem definition.** The target application is considered
 87 as a two-tuple $\mathcal{A} = (\mathcal{T}, \mathcal{D})$. Here, $\mathcal{T} = \{t_1, t_2, \dots, t_T\}$ denotes a set of $|\mathcal{T}| = T$ inde-
 88 pendent computational tasks and $\mathcal{D} = \{d_1, d_2, \dots, d_D\}$ denotes a set of $|\mathcal{D}| = D$ data
 89 elements. Here and hereafter $|\cdot|$ denotes the cardinality of the respective set. There
 90 exists no computational dependency between tasks. However, there exists interaction
 91 among the tasks as multiple tasks may need the same data element(s) for their exe-
 92 cutions. In a dual manner, an individual data element may be required by multiple

93 tasks for execution. The set of data elements required by a task t_i is denoted by
 94 $Data(t_i)$, whereas the set of tasks that need/require a data element d_j is denoted by
 95 $Tasks(d_j)$. Tasks may be associated with different computational costs, as well as
 96 data elements may be associated with different memory sizes. Let $exec(t_i)$ denote the
 97 computational cost of task t_i and let $size(d_j)$ denote the memory size of data element
 98 d_j .

99 Row-row-parallel Sparse Generalized Matrix Matrix Multiplication (SpGEMM)
 100 of the form $C=AB$ is an example target application. That is, the pre-multiplication
 101 of individual A -matrix rows with the B -matrix constitute the tasks, whereas rows of
 102 the A and B matrices constitute the data elements. The details about the SpGEMM
 103 application is given in subsection 5.3.2. Figure 1a shows a sample SpGEMM instance
 104 with a 3×4 A -matrix and a 4×5 B -matrix.

105 Figure 1b shows the $\mathcal{A}=(\mathcal{T}, \mathcal{D})$ representation of the SpGEMM instance given in
 106 Figure 1a with $|\mathcal{T}|=3$ tasks and $|\mathcal{D}|=7$ data elements. In the figure, circles denote
 107 tasks and squares denote data elements, whereas lines denote the interaction among
 108 the tasks and data elements. Vertices d_1, d_2 , and d_3 respectively denote A -matrix rows
 109 r_1^A, r_2^A , and r_3^A , whereas d_4, d_5, d_6 , and d_7 respectively denote B -matrix rows $r_1^B, r_2^B,$
 110 r_3^B , and r_4^B . The multiplication of the second row of the A -matrix with the B -matrix
 111 is represented by t_2 . This multiplication requires an A -matrix row r_2^A and three B -
 112 matrix rows r_2^B, r_3^B , and r_4^B . Therefore, $Data(t_2)=\{d_2, d_5, d_6, d_7\}$. The B -matrix row
 113 r_4^B is required by the first and second rows of the A -matrix, so $Tasks(d_7)=\{t_2, t_3\}$.
 114 $exec(t_i)$ and $size(d_j)$ values are also given for each element, for example $exec(t_2)=8$
 115 since it consists of $1+3+4=8$ multiply-add operations and $size(d_2)=3$ since r_2^A
 116 contains 3 nonzero elements.

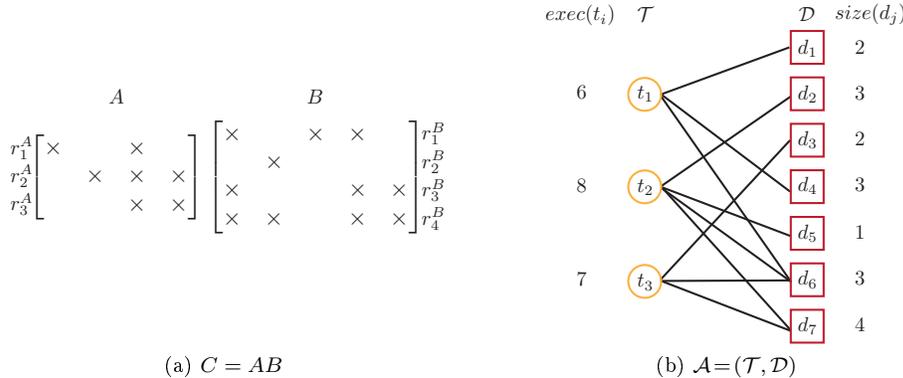


FIG. 1. A sample SpGEMM instance and corresponding $\mathcal{A}=(\mathcal{T}, \mathcal{D})$ representation.

117 The target computing platform is a homogeneous distributed-memory parallel
 118 system consisting of K processors. The execution time of each computational task is
 119 assumed to be the same on every processor.

120 The target problem is to find a computational-task-to-processor partition/assignment. ■
 121 Let $\Pi(\mathcal{T})=\{\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_K\}$ denote a K -way computational-task-to-processor assign-
 122 ment, where \mathcal{T}_k denotes the set of tasks assigned to processor p_k , for $k=1, 2, \dots, K$.
 123 This task-to-processor assignment $\Pi(\mathcal{T})$ incurs a data replication schema determined
 124 by the data needs of the tasks assigned to individual processors. In this schema, each
 125 data element d_j is replicated to each processor where at least one task assigned to
 126 that processor needs d_j . Let $WS(p_k)$ denote the working set of processor p_k , which
 127 corresponds to the set of data elements needed by p_k for the execution of the tasks

128 assigned to p_k . That is,

$$129 \quad (2.1) \quad WS(p_k) = \bigcup_{t_i \in \mathcal{T}_k} Data(t_i).$$

130 Consider a given task-to-processor assignment $\Pi(\mathcal{T})$. The computational load
131 $CL(p_k)$ of processor p_k is computed as

$$132 \quad (2.2) \quad CL(p_k) = \sum_{t_i \in \mathcal{T}_k} exec(t_i).$$

133 The data load $DL(p_k)$ of processor p_k is computed as

$$134 \quad (2.3) \quad DL(p_k) = \sum_{d_j \in WS(p_k)} size(d_j).$$

135 Note that $DL(p_k)$ corresponds to the memory footprint, which refers to the amount
136 of main memory that processor p_k references while executing the tasks assigned to
137 itself.

138 In a given assignment $\Pi(\mathcal{T})$, the maximum computational load and the maximum
139 data load of processors are respectively defined as

$$140 \quad (2.4) \quad CL_{max} = \max_k \{CL(p_k)\},$$

$$141 \quad (2.5) \quad DL_{max} = \max_k \{DL(p_k)\}.$$

143 After describing the framework and giving the notations, we formally define the
144 target problem as follows:

145 **DEFINITION 1** (Simultaneous Computation and Data Load Balancing Problem).

146 *Given an application $\mathcal{A} = (\mathcal{T}, \mathcal{D})$ find a computation-task-to-processor assignment*
147 *$\Pi(\mathcal{T}_k)$ that minimizes both CL_{max} and DL_{max} given in (2.4) and (2.5), respectively.*

148 3. Preliminaries.

149 **3.1. Hypergraph partitioning.** A hypergraph $\mathcal{H} = (\mathcal{U}, \mathcal{N})$ consists of a set \mathcal{U}
150 of vertices and a set \mathcal{N} of nets. Each net n_j connects a subset of vertices denoted as
151 $Pins(n_j)$. The degree $deg(n_j)$ of a net n_j denotes the number of vertices it connects,
152 i.e., $deg(n_j) = |Pins(n_j)|$. A cost $c(n_j)$ is associated with each net n_j . $Nets(u_i)$
153 denotes the set of nets that connect u_i . This easily extends to a subset of vertices
154 $\mathcal{U}_k \subset \mathcal{U}$ so that $Nets(\mathcal{U}_k) = \bigcup_{u_i \in \mathcal{U}_k} Nets(u_i)$. Multiple weights $w^1(u_i), \dots, w^C(u_i)$
155 can be associated with each vertex u_i , where $w^c(u_i)$ denotes the c th weight associated
156 with u_i .

157 $\Pi(\mathcal{H}) = \{\mathcal{U}_1, \dots, \mathcal{U}_K\}$ is called K -way partition of \mathcal{H} , if parts are mutually disjoint
158 and mutually exhaustive. In $\Pi(\mathcal{H})$, the connectivity set $\Lambda(n_j)$ of net n_j consists of
159 the parts that are connected by that net, i.e., $\Lambda(n_j) = \{\mathcal{U}_k : Pins(n_j) \cap \mathcal{U}_k \neq \emptyset\}$. The
160 number of parts connected by n_j is denoted by $\lambda(n_j) = |\Lambda(n_j)|$. A net n_j is said to be
161 cut if it connects more than one part, i.e., $\lambda(n_j) > 1$, and uncut otherwise. A vertex
162 u_i in $\Pi(\mathcal{H})$ is said to be a boundary vertex if it is connected by at least one cut net.
163 Among various cutsizes definitions we focus on the connectivity metric as follows:

$$164 \quad (3.1) \quad Cutsizes(\Pi(\mathcal{H})) = \sum_{n_j \in \mathcal{N}} c(n_j)(\lambda(n_j) - 1).$$

165 In a given partition $\Pi(\mathcal{H})$, the weight $W^c(\mathcal{U}_k)$ of part \mathcal{U}_k is defined as the sum of
 166 the c th weights of the vertices in \mathcal{U}_k . $\Pi(\mathcal{H})$ is said to be balanced if

$$167 \quad (3.2) \quad W^c(\mathcal{U}_k) \leq W_{avg}^c(1 + \epsilon^c), \text{ for } k \in \{1, 2, \dots, K\} \text{ and } c \in \{1, 2, \dots, C\},$$

168 where $W_{avg}^c = (\sum_k W^c(\mathcal{U}_k))/K$ and ϵ^c is the predetermined imbalance ratio for the
 169 c th weight.

170 The K -way multi-constraint HP problem [15] is then defined as finding a K -way
 171 partition such that the cutsizes (3.1) is minimized while the balance constraint (3.2) is
 172 maintained. For $C=1$, this reduces to the well-studied standard partitioning problem.

173 **3.2. Recursive bipartitioning (RB) framework.** In the RB paradigm, the
 174 initial hypergraph is partitioned into two subhypergraphs. These two subhypergraphs
 175 are further bipartitioned recursively until K parts are obtained. This process forms
 176 a complete binary tree, which we refer to as an RB tree, with $\log_2 K$ levels, where K
 177 is a power of 2.

178 The RB-based HP tools/algorithms utilize a cut-net splitting scheme in order to
 179 correctly encode the total cutsizes (3.1) at the end of the multi-way partitioning. That
 180 is, after each RB step, the cut nets of the respective vertex bipartition $\Pi_2 = \{\mathcal{U}_1, \mathcal{U}_2\}$
 181 are split into the two parts of the bipartition. Then, vertex-induced (induced by
 182 \mathcal{U}_1 and \mathcal{U}_2) subhypergraphs $\mathcal{H}_1 = (\mathcal{U}_1, \mathcal{N}_1)$ and $\mathcal{H}_2 = (\mathcal{U}_2, \mathcal{N}_2)$ of $\mathcal{H} = (\mathcal{U}, \mathcal{N})$ are
 183 constructed as follows:

$$184 \quad \mathcal{H}_1 = (\mathcal{U}_1, \mathcal{N}_1) \quad \mathcal{N}_1 = \{n' : \forall n \in \mathcal{N}, Pins(n) \cap \mathcal{U}_1 \neq \emptyset, Pins(n') = Pins(n) \cap \mathcal{U}_1\},$$

$$185 \quad \mathcal{H}_2 = (\mathcal{U}_2, \mathcal{N}_2) \quad \mathcal{N}_2 = \{n'' : \forall n \in \mathcal{N}, Pins(n) \cap \mathcal{U}_2 \neq \emptyset, Pins(n'') = Pins(n) \cap \mathcal{U}_2\}.$$

186 **4. Proposed hypergraph models.** In this section, we describe the two differ-
 187 ent hypergraph models proposed for solving the simultaneous computation and data
 188 load balancing problem.

189 **4.1. Hypergraph with data vertices (DV).** In this model, the application
 190 $\mathcal{A} = (\mathcal{T}, \mathcal{D})$ is represented by a hypergraph $\mathcal{H}_{DV}(\mathcal{A}) = (\mathcal{U} \cup \mathcal{V}, \mathcal{N})$ on $|\mathcal{T}| + |\mathcal{D}|$ vertices
 191 and $|\mathcal{D}|$ nets with the number of pins equal to

$$192 \quad (4.1) \quad \sum_{n \in \mathcal{N}} |Pins(n)| = \sum_{t_i \in \mathcal{T}} |Data(t_i)| + |\mathcal{D}| = \sum_{d_j \in \mathcal{D}} |Tasks(d_j)| + |\mathcal{D}|.$$

193 Vertex set \mathcal{U} represents the computational tasks, where each computational task t_i is
 194 represented by a task vertex $u_i \in \mathcal{U}$. Vertex set \mathcal{V} and net set \mathcal{N} represent the data
 195 elements. That is, each data element d_j is represented by a data vertex $v_j \in \mathcal{V}$ as well
 196 as a net $n_j \in \mathcal{N}$. Each net n_j connects the set of vertices representing the tasks that
 197 require data element d_j for their execution as well as data vertex v_j . That is,

$$198 \quad (4.2) \quad Pins(n_j) = \{v_j\} \cup \{u_i : t_i \in Tasks(d_j)\} = \{v_j\} \cup \{u_i : d_j \in Data(t_i)\}.$$

199 Therefore, each net connects one data vertex and one or more task vertices. Each
 200 net is associated with a cost which is equal to the memory size of the respective data
 201 element. That is,

$$202 \quad (4.3) \quad c(n_j) = size(d_j).$$

203 Without loss of generality, a given vertex partition $\Pi(\mathcal{H}_{DV}) = \{\mathcal{U}_1 \cup \mathcal{V}_1, \mathcal{U}_2 \cup$
 204 $\mathcal{V}_2, \dots, \mathcal{U}_K \cup \mathcal{V}_K\}$ is decoded as a K -way task assignment, where the tasks corre-
 205 sponding to the vertices in \mathcal{U}_k are assigned to processor p_k for $k = 1, \dots, K$. That

206 is $\mathcal{T}_k = \{t_i : u_i \in \mathcal{U}_k\}$. In this setting, $\Lambda(n_j)$ is interchangeably used for both the
 207 parts that net n_j connects and the respective processors. Recall that this vertex
 208 partition also incurs a data assignment/replication schema. A data element d_j is as-
 209 signed/replicated to each processor p_k such that $\mathcal{U}_k \in \Lambda(n_j)$. In other words, for each
 210 net n_j that connects part \mathcal{U}_k , data element d_j is assigned/replicated to processor p_k .
 211 That is,

$$212 \quad (4.4) \quad WS(p_k) = \{d_j : n_j \in Nets(\mathcal{U}_k)\}.$$

213 For a given partition $\Pi(\mathcal{H}_{DV})$, consider an internal net n_j in $\mathcal{U}_k \cup \mathcal{V}_k$. Then, all
 214 tasks which need data element d_j are assigned to the same processor p_k which already
 215 holds d_j . So, internal nets do not incur any replication.

216 Consider a cut net n_j and assume that v_j is assigned to $\mathcal{V}_k \in \Lambda(n_j)$. Then two
 217 cases occur as follows: Net n_j connects at least one task vertex in \mathcal{U}_k , net n_j does
 218 not connect any task vertex in \mathcal{U}_k . In the former case, each processor in $\Lambda(n_j)$ needs
 219 data element d_j , whereas in the latter case, each processor $\Lambda(n_j) \setminus \{p_k\}$ needs data
 220 element d_j . So in both cases, the data element d_j will be replicated to all processors
 221 in $\Lambda(n_j) \setminus \{p_k\}$. Hence, $c(n_j)(\lambda(n_j) - 1)$ denotes the total amount of replication
 222 because of the data element d_j . So, the partitioning objective of minimizing the
 223 cutsizes (3.1) corresponds to minimizing the total data replication to be incurred by
 224 the task partition.

225 Data vertices are included in the pin lists of the respective nets (e.g., $v_j \in$
 226 $Pins(n_j)$ as shown in (4.2)) for encoding the data loads of the processors through
 227 a two-constraint partitioning formulation as follows: The first weight of a task vertex
 228 is set equal to the execution time of the respective task, whereas its second weight is
 229 set to zero. The first weight of a data vertex is set to zero, whereas its second weight
 230 is set equal to the memory size of the respective data element. That is,

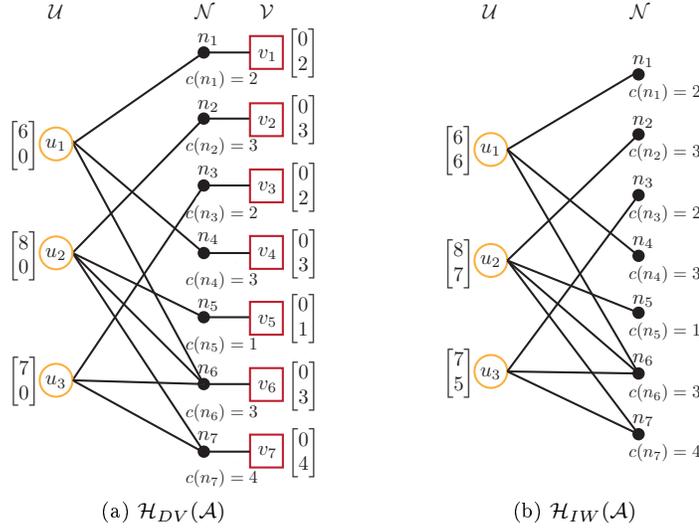
$$231 \quad (4.5) \quad w^1(u_i) = exec(t_i) \quad w^1(v_j) = 0$$

$$232 \quad (4.6) \quad w^2(u_i) = 0 \quad w^2(v_j) = size(d_j).$$

233 So, the first partitioning constraint of maintaining balance on parts' first weights
 234 encodes balancing computational loads of processors, whereas the second partitioning
 235 constraint of maintaining balance on parts' second weights relates to balancing data
 236 loads of processors.

237 Figure 2a shows the \mathcal{H}_{DV} hypergraph for the sample application $\mathcal{A} = (\mathcal{T}, \mathcal{D})$
 238 given in Figure 1. In the figure, circles denote task vertices, squares denote data
 239 vertices and dots denote the nets, whereas lines denote pins. For example, $Pins(n_6) =$
 240 $\{u_1, u_2, u_3, v_6\}$ since $Tasks(d_6) = \{t_1, t_2, t_3\}$. The array of two weights associated with
 241 each vertex is displayed next to the corresponding vertex, where the upper weight
 242 denotes $w^1(\cdot)$ and the lower weight denotes $w^2(\cdot)$. For example, $w^1(u_2) = 8$ since
 243 $exec(t_2) = 8$ and $w^2(v_6) = 3$ since $size(d_6) = 3$.

244 In a straightforward partitioning of \mathcal{H}_{DV} , the relation between second constraint
 245 and balancing processors' data loads is rather loose since each data vertex is assigned
 246 to only one part and does not encode replication of data elements according to the
 247 task partition. We enhance this two-constraint formulation with a novel boundary
 248 data vertex replication scheme utilized in an RB framework for enabling the second
 249 constraint to better encode balancing processors' data loads. In this scheme, the
 250 bipartition of the computational task vertices obtained at each RB step is utilized to
 251 determine data vertex replication in the further RB steps. The data vertex replication
 252 is performed together with the conventional cut-net splitting scheme as follows:

FIG. 2. Hypergraph models for the sample application $\mathcal{A}=(\mathcal{T}, \mathcal{D})$ given in Figure 1.

253 Consider a bipartition $\Pi_2 = \{\mathcal{U}_1 \cup \mathcal{V}_1, \mathcal{U}_2 \cup \mathcal{V}_2\}$ of \mathcal{H}_{DV} at the end of the current
 254 RB step. Consider a cut net n_j in Π_2 . Cut net n_j possibly connects task vertices
 255 in both parts, whereas it connects the respective data vertex v_j which is a boundary
 256 vertex in one of the two parts. In the conventional net splitting, net n_j will be split
 257 into both parts as n'_j and n''_j with $Pins(n'_j) = Pins(n_j) \cap (\mathcal{U}_1 \cup \mathcal{V}_1)$ and $Pins(n''_j) =$
 258 $Pins(n_j) \cap (\mathcal{U}_2 \cup \mathcal{V}_2)$, respectively. In the proposed scheme, boundary vertex v_j in one
 259 part will be replicated to the other part (as v'_j) so that both of the split nets n'_j and
 260 n''_j connect data vertex v_j or v'_j , both of which represent data element d_j . That is,

$$261 \quad (4.7) \quad Pins(n'_j) = (Pins(n_j) \cap \mathcal{U}_1) \cup \{v_j\},$$

$$262 \quad (4.8) \quad Pins(n''_j) = (Pins(n_j) \cap \mathcal{U}_2) \cup \{v'_j\}.$$

263 So, bipartition $\Pi_2 = \{\mathcal{U}_1 \cup \mathcal{V}_1, \mathcal{U}_2 \cup \mathcal{V}_2\}$ obtained at a particular RB step induces the
 264 hypergraphs \mathcal{H}_{DV1} and \mathcal{H}_{DV2} for further bipartitioning in the following RB steps:

$$265 \quad \mathcal{H}_{DV1} = ((\mathcal{U}_1 \cup \mathcal{V}_1) \cup \mathcal{V}_2^B, \mathcal{N}_1) \text{ and } \mathcal{H}_{DV2} = ((\mathcal{U}_2 \cup \mathcal{V}_2) \cup \mathcal{V}_1^B, \mathcal{N}_2), \text{ where}$$

$$266 \quad \mathcal{N}_1 = \{n' : \forall n \in \mathcal{N}, Pins(n) \cap (\mathcal{U}_1 \cup \mathcal{V}_1) \neq \emptyset, Pins(n') = Pins(n) \cap (\mathcal{U}_1 \cup \mathcal{V}_1 \cup \mathcal{V}_2^B)\}$$

$$267 \quad \mathcal{N}_2 = \{n'' : \forall n \in \mathcal{N}, Pins(n) \cap (\mathcal{U}_2 \cup \mathcal{V}_2) \neq \emptyset, Pins(n'') = Pins(n) \cap (\mathcal{U}_2 \cup \mathcal{V}_2 \cup \mathcal{V}_1^B)\}.$$

268 Here, \mathcal{V}_1^B and \mathcal{V}_2^B respectively denote boundary data vertex sets of \mathcal{V}_1 and \mathcal{V}_2 . So,
 269 $((\mathcal{U}_1 \cup \mathcal{V}_1) \cup \mathcal{V}_2^B)$ denotes the replication of boundary data vertex set of \mathcal{V}_2 to \mathcal{H}_{DV1} and
 270 in a dual manner $(\mathcal{U}_2 \cup \mathcal{V}_2) \cup \mathcal{V}_1^B$ denotes the replication of the boundary data vertex
 271 set of \mathcal{V}_1 to \mathcal{H}_{DV2} . The weights of vertices remain the same after the RB step.

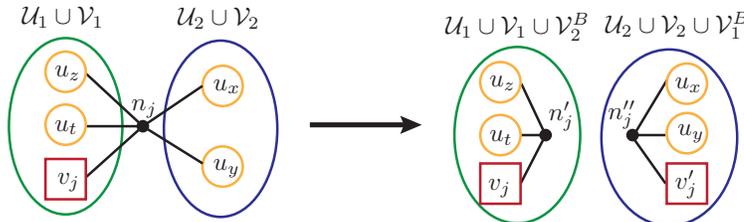


FIG. 3. Proposed data vertex replication together with cut-net splitting.

272 **Figure 3** shows the usage of the RB framework through a sample bipartition in
 273 terms of a single cut net n_j which connects two task vertices and one data vertex in
 274 the left part and two task vertices in the right part. Here and hereafter part 1 and
 275 part 2 of the bipartition are respectively referred to as left and right parts for clarity
 276 of the presentation. As seen in the figure, this cut net is split as n' and n'' to the
 277 left and right parts, respectively. The boundary data vertex v_j connected by n_j in
 278 the left part is replicated to the right part as v'_j so that $Pins(n'_j) = \{u_x, u_t, v_j\}$ and
 279 $Pins(n''_j) = \{u_x, u_y, v'_j\}$.

280 In the conventional cut-net splitting scheme [14], after an RB step, if a cut net
 281 n_j connects only one vertex in one of the parts, then n_j is not split to that part since
 282 it will incur a single-pin net in that part and single-pin nets do not contribute to the
 283 cutsizes in the further RB steps. However, in the proposed scheme, such cases should
 284 be handled differently depending on whether the only vertex connected by a cut net
 285 in one of the parts of the bipartition is a data vertex or task vertex.

286 **Special case 1:** This case occurs when the only vertex connected by a cut net
 287 n_j in one part is a data vertex v_j . The proposed scheme replicates v_j to the other
 288 part. However, vertex v_j connected by a single-pin split net n'_j corresponds to a data
 289 element assigned to a processor that is not assigned any computation task which needs
 290 the data element d_j . Hence, we move (instead of replicating) data vertex v_j to the
 291 other part so that the cut net n_j becomes internal. Note that this move operation will
 292 decrease the cutsizes but has the potential of increasing the imbalance on the second
 293 part weights.

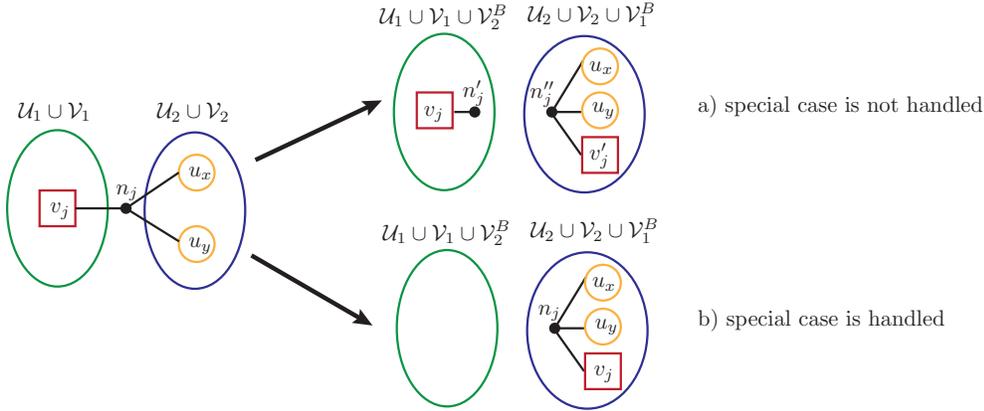


FIG. 4. *Special case 1: the only vertex connected by a cut net in one of the parts of the bipartition is a data vertex. a) special case not handled, b) special case handled*

294 **Figure 4** shows the handling of the special case 1 through a sample bipartition
 295 in the context of a single cut net n_j . Cut net n_j connects only a data vertex v_j on
 296 the left part, whereas it connects two task vertices u_x and u_y on the right part. The
 297 upper arrow shows the splitting of n_j as well as the replication of v_j , if the special
 298 case is not handled. The lower arrow shows the proposed handling of the special case,
 299 where the data vertex v_j on the left part moved to the right part so that n_j becomes
 300 an internal net of the right part.

301 **Special case 2:** This case occurs when the only vertex connected by a cut net
 302 n_j in one part is a task vertex u_i . The proposed scheme replicates v_j in the other
 303 part to this part. However, this replication will incur a two-pin split net connecting
 304 task vertex u_i and data vertex v'_j , which refers to a data vertex needed by a single
 305 task vertex. The trivial solution for such two-pin nets is to maintain them inter-

306 nal in further RB steps by avoiding this net splitting together with the data vertex
 307 replication while assigning the second weight of data vertices to the second weight of
 308 the task vertices. That is, u_i will contain two nonzero weights $w^1(u_i) = exec(t_i)$ and
 309 $w^2(u_i) = w^2(u_i) + size(d_j)$, on the contrary the initial two-constraint formulation in
 310 the top-most level where each vertex has one nonzero and one zero weight.

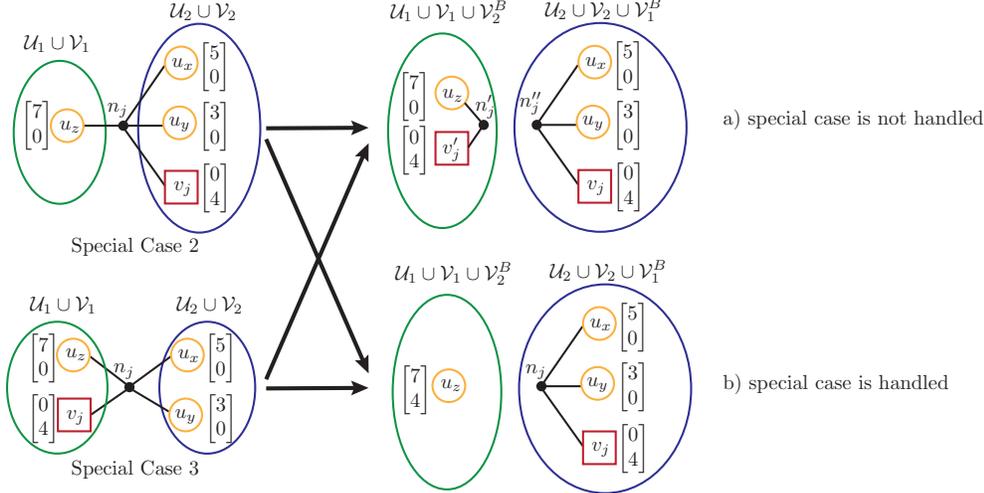


FIG. 5. *Special case 2: the only vertex connected by a cut net in one of the parts of the bipartition is a task vertex, whereas special case 3: a cut net connects only one task vertex together with the data vertex in one of the parts. a) special case not handled, b) special case handled*

311 **Figure 5** shows the handling of special case 2 after a sample RB step through a
 312 sample bipartition in the context of a single cut net n_j . Cut net n_j connects only a
 313 task vertex on the left part, whereas it connects two task vertices and a data vertex
 314 on the right part. The upper arrow shows the splitting of n_j as well as the replication
 315 of v_j , if the special case is not handled. The lower arrow shows the proposed handling
 316 of the special case, where neither v_j is replicated to the left part nor n_j is split to the
 317 left part. Instead the second weight of data vertex v_j on the right part is added to
 318 the second weight of the task vertex on the left part.

319 **Special case 3:** This case occurs when a cut net n_j connects only one task vertex
 320 u_z together with the data vertex v_j in one of the parts. The proposed scheme will
 321 incur a two-pin split net connecting task vertex u_z and data vertex v_j in that part. So
 322 this case becomes very similar to the special case 2 and handled in the same manner
 323 by moving v_j to the other part and add its second weight to the second weight of u_z .

324 **Figure 5** shows the handling of special case 3 which is equivalent to the special
 325 case 2 except moving v_j from the right part to the left part so that handled and
 326 unhandled split partitions will be the same.

327 **Algorithm 4.1** shows the RB-based partitioning of \mathcal{H}_{DV} utilizing the proposed
 328 data vertex replication scheme. If-statements at lines 11-14, 15-19, and 20-26 respec-
 329 tively show the handling of the special cases 1, 2, and 3. Statements at lines 28-30
 330 show the replication of the data vertex to the other part. Note that data vertex
 331 replication is performed only if none of the special cases occur for the respective net.

332 **4.2. Hypergraph model with inverse data weight (IW) distribution.** In
 333 this model, the application $\mathcal{A} = (\mathcal{T}, \mathcal{D})$ is represented by a hypergraph $\mathcal{H}_{IW}(\mathcal{A}) = (\mathcal{U}, \mathcal{N})$ ■

Algorithm 4.1 Partition \mathcal{H}_{DV} with proposed data vertex replicationInput: $\mathcal{H}_{DV} = (\mathcal{U} \cup \mathcal{V}, \mathcal{N}, w^1, w^2)$, K Output: $\Pi(\mathcal{H}_{DV})$

```

1:  $\mathcal{H}_0^0 = \mathcal{H}_{DV}$ 
2: for  $\ell \leftarrow 0$  to  $\log_2 K - 1$  do
3:   for  $k \leftarrow 0$  to  $2^\ell - 1$  do
4:      $\Pi_2 \leftarrow \text{BIPARTITION}(\mathcal{H}_k^\ell)$  ▷  $\Pi_2 = \{\mathcal{U}_L \cup \mathcal{V}_L, \mathcal{U}_R \cup \mathcal{V}_R\}$ 
5:     for each cut net  $n_j \in \mathcal{N}_k^\ell$  do
6:        $flag \leftarrow \text{true}$ 
7:       if  $v_j \in \mathcal{V}_L$  then
8:          $x \leftarrow L; y \leftarrow R$ 
9:       else
10:         $x \leftarrow R; y \leftarrow L$ 
11:       if  $Pins(n_j) \cap \mathcal{U}_x = \emptyset$  then ▷ special case 1
12:          $\mathcal{V}_y \leftarrow \mathcal{V}_y \cup \{v_j\}$ 
13:          $\mathcal{V}_x \leftarrow \mathcal{V}_x \setminus \{v_j\}$  ▷ move data vertex  $v_j$  to the other part
14:          $flag \leftarrow \text{false}$  ▷ Net  $n_j$  becomes internal
15:       if  $|Pins(n_j) \cap \mathcal{U}_y| = 1$  then ▷ special case 2
16:          $u_i \leftarrow Pins(n_j) \cap \mathcal{U}_y$ 
17:          $w^2(u_i) \leftarrow w^2(u_i) + w^2(v_j)$ 
18:          $Pins(n_j) \leftarrow Pins(n_j) \setminus \{u_i\}$ 
19:          $flag \leftarrow \text{false}$  ▷ Net  $n_j$  becomes internal
20:       if  $|Pins(n_j) \cap \mathcal{U}_x| = 1$  then ▷ special case 3
21:          $u_i \leftarrow Pins(n_j) \cap \mathcal{U}_x$ 
22:          $w^2(u_i) \leftarrow w^2(u_i) + w^2(v_j)$ 
23:          $\mathcal{V}_y \leftarrow \mathcal{V}_y \cup \{v_j\}$ 
24:          $\mathcal{V}_x \leftarrow \mathcal{V}_x \setminus \{v_j\}$  ▷ move data vertex  $v_j$  to the other part
25:          $Pins(n_j) \leftarrow Pins(n_j) \setminus \{u_i\}$ 
26:          $flag \leftarrow \text{false}$  ▷ Net  $n_j$  becomes internal
27:       if  $flag$  then
28:          $v'_j \leftarrow v_j$ 
29:          $\mathcal{V}_y \leftarrow \mathcal{V}_y \cup \{v'_j\}$  ▷ replicate  $v_j$  to other part
30:          $Pins(n_j) \leftarrow Pins(n_j) \cup \{v'_j\}$ 
31:       Form  $\mathcal{H}_{2k}^{\ell+1} = (\mathcal{U}_L \cup \mathcal{V}_L, \mathcal{N}_L)$  induced by  $\mathcal{U}_L \cup \mathcal{V}_L$ 
32:       Form  $\mathcal{H}_{2k+1}^{\ell+1} = (\mathcal{U}_R \cup \mathcal{V}_R, \mathcal{N}_R)$  induced by  $\mathcal{U}_R \cup \mathcal{V}_R$ 

```

334 on $|\mathcal{T}|$ vertices and $|\mathcal{D}|$ nets with the number of pins equals to

$$335 \quad (4.9) \quad \sum_{n \in \mathcal{N}} |Pins(n)| = \sum_{t_i \in \mathcal{T}} |Data(t_i)| = \sum_{d_j \in \mathcal{D}} |Tasks(d_j)|.$$

336 Vertex set \mathcal{U} represents the computational tasks. Net set \mathcal{N} represents the data
 337 elements. That is, each computational task t_i is represented by a task vertex $u_i \in \mathcal{U}$
 338 and each data element d_j is represented by a net $n_j \in \mathcal{N}$. Each net n_j connects
 339 the set of vertices representing the tasks that require the data element d_j for their
 340 execution. That is,

$$341 \quad (4.10) \quad Pins(n_j) = \{u_i : t_i \in Tasks(d_j)\} = \{u_i : d_j \in Data(t_i)\}.$$

342 Comparison of (4.10) and (4.2) shows that \mathcal{H}_{DV} and \mathcal{H}_{IW} topologically differs

343 by vertex set \mathcal{V} , which corresponds to data elements of \mathcal{H}_{DV} . That is, \mathcal{H}_{DV} becomes
 344 topologically the same with the \mathcal{H}_{IW} when data vertices in \mathcal{V} and corresponding pins
 345 are removed.

346 Each net is associated with a cost which is equal to the memory size of the
 347 respective data element. That is,

$$348 \quad (4.11) \quad c(n_j) = size(d_j).$$

349 Without loss of generality, a given vertex partition $\Pi(\mathcal{H}_{IW}) = \{\mathcal{U}_1, \mathcal{U}_2, \dots, \mathcal{U}_K\}$ is
 350 decoded as a K -way task partition, where the tasks corresponding to the vertices in
 351 \mathcal{U}_k are assigned to processor p_k for $k=1, \dots, K$. Recall that this vertex partition also
 352 incurs a data assignment/replication schema. In other words, for each net n_j that
 353 connects part \mathcal{U}_k , data element d_j is assigned/replicated to processor p_k . That is,

$$354 \quad (4.12) \quad WS(p_k) = \{d_j : n_j \in Nets(\mathcal{U}_k)\}.$$

355 For a given partition $\Pi(\mathcal{H}_{IW})$, $\lambda(n_j)$ denotes the number of processors that need
 356 the data element d_j . So $\lambda(n_j) - 1$ denotes the number of times the data element d_j
 357 needs to be replicated. Hence, $c(n_j)(\lambda(n_j) - 1)$ denotes the total amount of replication
 358 because of the data element d_j . So, the partitioning objective of minimizing the cutsizes
 359 according to (3.1) corresponds to minimizing the total data replication to be incurred
 360 by the task partition.

361 In this model, a two-constraint partitioning formulation is also used, where the
 362 first and second weights of each vertex refer to the computational and data loads of
 363 the respective task. Since \mathcal{H}_{IW} does not contain data vertices, we propose a novel
 364 inverse data weight distribution model for estimating the second weights of vertices.
 365 In this model, the cost of a net, which corresponds to the size of the respective data
 366 element, is distributed evenly among the second weights of the vertices connected by
 367 that net. That is, a net n_j of cost $c(n_j)$, which represents data element d_j of $size(d_j)$,
 368 contributes $c(n_j)/deg(n_j)$ to each vertex it connects. Finally, the first weight of a task
 369 vertex is set equal to the execution time of the respective task, whereas its second
 370 weight is set to sum of the contributions from each net n_j connecting that vertex by
 371 $c(n_j)/deg(n_j)$. That is,

$$372 \quad (4.13) \quad w^1(u_i) = exec(t_i) \quad w^2(u_i) = \sum_{n_j \in Nets(u_i)} \frac{c(n_j)}{deg(n_j)}$$

374 So, the first partitioning constraint of maintaining balance on parts' first weights
 375 encodes balancing computational loads of processors, whereas the second partitioning
 376 constraint of maintaining balance on parts' second weights relates to balancing data
 377 loads of processors.

378 **Figure 2b** shows the \mathcal{H}_{IW} hypergraph for the sample application $\mathcal{A} = (\mathcal{T}, \mathcal{D})$
 379 given in **Figure 1**. For example, $Pins(n_7) = \{u_2, u_3\}$ since $Tasks(d_7) = \{t_2, t_3\}$. For
 380 example, $w^1(u_2) = 8$ since $exec(t_2) = 8$. Regarding the second weight of u_2 ; the nets
 381 n_2, n_5, n_6 , and n_7 , which connect u_2 , respectively contribute $c(n_2)/deg(n_2) = 3/1 = 3$,
 382 $c(n_5)/deg(n_5) = 1/1 = 1$, $c(n_6)/deg(n_6) = 3/3 = 1$, and $c(n_7)/deg(n_7) = 4/2 = 2$ to
 383 $w^2(u_2)$. That is, $w^2(u_2) = 3 + 1 + 1 + 2 = 7$.

384 The motivation behind the proposed inverse data weight distribution model can
 385 be described as follows: In a given partition of \mathcal{H}_{IW} , consider an internal net n_j of
 386 \mathcal{U}_k . This net refers to the case where all tasks requiring data element d_j are assigned
 387 to the same processor p_k . Net n_j will contribute a total weight of $c(n_j)$ to the second

388 weight of \mathcal{U}_k which in turn will correspond to contributing $c(n_j) = \text{size}(d_j)$ to the
 389 data load $DL(p_k)$ of processor p_k . So, internal nets correctly encode the data loads
 390 of the processors to which they are internal.

391 However, consider a cut net n_j with connectivity set $\Lambda(n_j)$. Net n_j will contribute
 392 fractional weights to the second weights of parts/processors in $\Lambda(n_j)$. The distribution
 393 of the weight $c(n_j)$ will be proportional to the number of pins it connects in those
 394 parts. That is, for each $\mathcal{U}_k \in \Lambda(n_j)$, net n_j will contribute

$$395 \quad (4.14) \quad \frac{|\text{Pins}(n_j) \cap \mathcal{U}_k|}{|\text{Pins}(n_j)|} c(n_j)$$

396 to the second weight $W^2(\mathcal{U}_k)$ of part \mathcal{U}_k . This in turn corresponds to net n_j con-
 397 tributing $\text{size}(d_j)|\text{Pins}(n_j) \cap \mathcal{U}_k|/\text{deg}(n_j)$ to the data load $DL(p_k)$ of processor p_k .

398 That is, data load of a processor is correctly encoded by internal nets in the
 399 corresponding part, whereas an error is made by the cut nets connecting that part.
 400 This is because data weight is encoded partially by the vertices assigned to that part.
 401 Note that the partitioning objective of minimizing the cutsize will minimize this error
 402 due to the cut nets. Also, errors made due to the fractional weight distribution of the
 403 cut nets can be expected to cancel each other. Consider two nets n_j and n_h of equal
 404 degree $\text{deg}(n_j) = \text{deg}(n_h) = \mathbf{deg}$ and equal cost $c(n_j) = c(n_h) = \mathbf{size}$. Assume that in
 405 the given partition, these two nets become cut and connect only the same two parts
 406 \mathcal{U}_k and \mathcal{U}_ℓ in the partition. Also assume that n_j connects α pins in \mathcal{U}_k and $\mathbf{deg} - \alpha$ pins
 407 in \mathcal{U}_ℓ , whereas n_h connects $\mathbf{deg} - \alpha$ pins in \mathcal{U}_k and α pins in \mathcal{U}_ℓ . Despite the erroneous
 408 fractional data load contributions because of these two cut nets to the data loads
 409 of processors p_k and p_ℓ , they together contribute the same amount of data load of
 410 \mathbf{size} to both processors p_k and p_ℓ . Although the actual aggregate contribution of n_j
 411 and n_h should be $2\mathbf{size}$ to both processors, assigning the same load of \mathbf{size} to both
 412 parts enables partitioner's load balancing mechanism to indirectly encode balancing
 413 data loads of processors. This discussion can be extended for the nets with different
 414 number of pins and costs.

415 Here we exploit the RB framework in order to improve the proposed \mathcal{H}_{IW} model as
 416 follows: Recall that the proposed model distributes the cost of each net evenly among
 417 the second weights of the vertices that it connects. The degrees of the nets decrease
 418 each time they become cut during the RB process because of the cut-net splitting
 419 scheme adopted. So, updated degree information of the split nets should be used
 420 for a more accurate net cost distribution. Therefore, after each RB step, the second
 421 weights of vertices in each of the two subhypergraphs are computed from scratch by
 422 taking into account the updated degree information of the split nets. Although the
 423 contributions of the internal nets to the second vertex weights do not change since
 424 their degrees remain the same, computing the second weights from scratch seems to
 425 be more efficient.

426 Figure 6 shows the usage of the RB framework through a sample bipartition in
 427 terms of a single cut net n_j of degree five. As seen in the figure, n_j evenly distributes
 428 its cost 30 among the second weights of those five vertices as $30/5 = 6$ before the
 429 current bipartitioning step. After the RB step, the degrees of the split nets n' and
 430 n'' become three and two, respectively. Therefore, in the left part, n'_j contributes
 431 $30/3 = 10$ to the second weights of the vertices u_z , u_t and, u_s , whereas, in the right
 432 part, n''_j contributes $30/2 = 15$ to the second weights of the vertices u_x and u_y .

433 Algorithm 4.2 shows the RB-based partitioning of \mathcal{H}_{IW} with the proposed inverse
 434 data weight distribution. The for loop at lines 6-9 computes the inverse data weight
 435 of each net and then distributes this weight to those vertices that it connects.

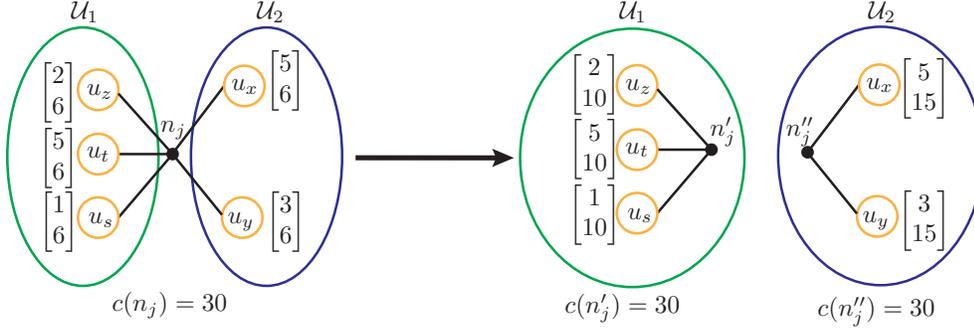


FIG. 6. Proposed inverse data weight distribution together with cut-net splitting.

Algorithm 4.2 Inverse Data Weight Distribution AlgorithmInput: $\mathcal{H}_{IW} = (\mathcal{U}, \mathcal{N}, w^1, c), K$ Output: $\Pi(\mathcal{H}_{IW})$

- 1: $\mathcal{H}_0^0 = \mathcal{H}_{IW}$
- 2: **for** $\ell \leftarrow 0$ **to** $\log_2 K - 1$ **do**
- 3: **for** $k \leftarrow 0$ **to** $2^\ell - 1$ **do**
- 4: **for** each vertex $u_i \in \mathcal{U}_k^\ell$ **do**
- 5: $w^2(u_i) \leftarrow 0$
- 6: **for** each net $n_j \in \mathcal{N}_k^\ell$ **do**
- 7: $idwContr \leftarrow c(n_j)/deg(n_j)$;
- 8: **for** each $u_i \in Pins(n_j)$ **do**
- 9: $w^2(u_i) \leftarrow w^2(u_i) + idwContr$
- 10: $\Pi_2 \leftarrow \text{BIPARTITION}(\mathcal{H}_k^\ell)$ $\triangleright \Pi_2 = \{\mathcal{V}_L, \mathcal{V}_R\}$
- 11: Form $\mathcal{H}_L = \mathcal{H}_{2k}^{\ell+1} = (\mathcal{V}_L, \mathcal{N}_L)$ induced by \mathcal{V}_L
- 12: Form $\mathcal{H}_R = \mathcal{H}_{2k+1}^{\ell+1} = (\mathcal{V}_R, \mathcal{N}_R)$ induced by \mathcal{V}_R

436 **4.3. Discussion.** As mentioned earlier, the two proposed models are topologi-
437 cally similar, where nets represent data elements. So, in both models, the partitioning
438 objective of minimizing the cutsizes encodes the minimization of the total amount of
439 data replication via clustering the tasks that require the same data elements to the
440 same parts under the given balancing constraints. This partitioning objective also
441 encodes the amount of communication volume to incur for realizing the required data
442 replication among processors. We should note here that the same data element re-
443 quired by multiple tasks assigned to the same processor necessitates the replication
444 of that data element only once. Both proposed hypergraph models encapsulate this
445 replication correctly, whereas a similar bipartite graph model would overestimate.

446 As discussed earlier, both models utilize a two-constraint formulation, where the
447 first and second weights of vertices are respectively used to encode the computational
448 and the data loads of parts/processors. Since the total amount of computation is
449 constant, the partitioning constraint on maintaining balance on the parts' first weights
450 encodes minimizing the computational load of the maximally loaded processor within
451 the given computational load imbalance ratio (ϵ^1). However, the total amount of data
452 load of processors, which is the sum of the second weights of the parts, is not constant
453 and it depends on the quality of the task partitioning. So, a naive balancing on the
454 seconds weights of the parts might not encode the minimization of the data load of
455 the maximally loaded processor. For example, a very tight balance on the data loads

456 of the processors might yield a very high data load on the maximally loaded processor
 457 if the underlying partition produces a high amount of data replication. Therefore,
 458 the second partitioning constraint on maintaining balance on parts' second weights
 459 under the partitioning objective that encodes the minimization of the total amount
 460 of data replication corresponds to minimizing the data load of the maximally loaded
 461 processor within the given data load imbalance ratio (ϵ^2).

462 The two hypergraph models differ in the second vertex weighting scheme they
 463 utilize for estimating processors' data loads. Both models correctly encode the pro-
 464 cessors' data loads corresponding to the data elements that are required by a single
 465 processor. Both models utilize the RB framework to increase the accuracy of their
 466 schemes utilized to estimate processors' data loads corresponding to the data elements
 467 that are required by multiple processors. After each RB step, the \mathcal{H}_{DV} model corrects
 468 the topology of the subhypergraphs by augmenting the subhypergraphs with the repli-
 469 cated data vertices according to the induced task vertex bipartition. At the beginning
 470 of each RB step, the \mathcal{H}_{IW} model utilizes the topology of the current hypergraph to
 471 predict the difference between data loads of the two parts of the bipartition. That is,
 472 the \mathcal{H}_{DV} model tries to encode data loads of processors after the RB step, whereas
 473 the \mathcal{H}_{IW} model tries to encode data loads of processors before the RB step.

474 5. Experiments.

475 **5.1. Baseline model.** The baseline model is the conventional HP model widely
 476 used for the parallelization of irregularly sparse applications. The topology of this
 477 hypergraph model is exactly the same with that of the \mathcal{H}_{IW} model described in [sub-](#)
 478 [section 4.2](#). However, this model is a single constraint model, where the vertices are
 479 weighted with the computational loads of the respective tasks they represent. In this
 480 way, the partitioning constraint of balancing part weights encodes computational bal-
 481 ance among processors. The cost of the nets is set to be equal to the size of the
 482 data elements they represent. So, the partitioning objective of minimizing the cuts-
 483 ize according to connectivity metric [\(3.1\)](#) encodes minimizing total communication vol-
 484 ume [\[15\]](#). Note that this partitioning objective also encodes minimizing total amount
 485 of data replication. So, the baseline model differs from the proposed models in not
 486 considering data load balancing at all. Here and hereafter, we refer to this hypergraph
 487 model as \mathcal{H}_{Base} .

488 **5.2. Experimental setup.** The hypergraph models proposed in [subsections 4.1](#)
 489 [and 4.2](#), as well as the baseline hypergraph model mentioned in [subsection 5.1](#) are par-
 490 titioned using the HP tool PaToH [\[14, 16\]](#) for obtaining $K \in \{64, 128, 256, 512, 1024\}$ -
 491 way partitions. PaToH is used with default parameters for all models except for
 492 the \mathcal{H}_{DV} model described in [subsection 4.1](#). For the \mathcal{H}_{DV} model, we set PaToH's
 493 vertex visit order to the continuous/sequential vertex order (increasing vertex ID or-
 494 der) for the coarsening phase instead of the random vertex visit order which is the
 495 default [\[16\]](#). The objective behind this is to prioritize matching of computational ver-
 496 tices with other computational or data vertices. In order to maintain randomness in
 497 vertex visit order, we randomly permute computational and data vertices separately
 498 before invoking PaToH. In all partitioning instances, we used maximum allowable
 499 imbalance ratio $\epsilon = 0.05$ for both computation and data weights, i.e., $\epsilon^1 = 0.05$ and
 500 $\epsilon^2 = 0.05$. All experiments have been conducted by use of random seed. As PaToH
 501 utilizes randomized algorithms, we partitioned each instance five times with different
 502 seeds and we report the geometric average of the results.

503 **5.3. Dataset.** The performance of the two proposed models are validated against
 504 the baseline model on two sample applications: Parallel Finite Element Method
 505 (FEM) and Volume Element Method (VEM) based simulations which involve par-
 506 titioning irregular 2D or 3D meshes and parallel Sparse Generalized Matrix Matrix
 507 Multiplication (SpGEMM) which involves partitioning two irregularly sparse input
 508 matrices.

509 **5.3.1. Mesh partitioning instances.** In the FEM/VEM applications, compu-
 510 tations associated with each mesh/volume element (cell) constitute an atomic com-
 511 putational task. These applications implement ghost layering methods which involve
 512 replicating cells according to the cell-to-cell neighborhood relation determined by the
 513 target application [22]. Hence, such FEM/VEM applications fall within the frame-
 514 work in section 2 since computations associated with each element are independent
 515 but they share data elements determined by the neighborhood relation.

516 In order to obtain mesh partitioning instances, we utilize sparse matrices which
 517 have either 2D or 3D coordinate values so that the sparsity patterns of those matrices
 518 are considered as representing the neighborhood structures of the meshes arising in
 519 FEM or VEM applications. Such matrices are selected from the SuiteSparse Matrix
 520 Collection [23]. These matrices are symmetric matrices so that rows and columns
 521 respectively represent mesh elements and data elements or vice versa. So, for a selected
 522 sparse matrix $A=(a_{ij})$, we have

$$523 \quad (5.1) \quad \text{Data}(t_i) = \{d_j : a_{ij} \neq 0\}.$$

524 Here, t_i denotes the atomic task associated with mesh cell c_i and d_j denotes the data
 525 associated with cell c_j .

526 We should note that SuiteSparse Matrix Collection [23] does not contain any
 527 information about computational cost and data size distribution. Various computa-
 528 tional cost and data size weighting schemes are utilized depending on the applica-
 529 tion nature of mesh computations [7,22,29,34]. In this paper, we use the weighting scheme
 530 produced by a heuristic for generating realistic weight distributions for “Particles-in-
 531 Cells”-like simulations [7]. The computational cost of a mesh element is reported to
 532 be equal to the square of its memory size for the weighting [34] and the amount of
 533 memory needed for holding a cell is linear with the number of particles located in
 534 this cell [7]. Reasoning behind this is data size is related to the number of particles
 535 located in a cell, while computational cost associated with a cell generally increases
 536 with the square of the number of particles in this cell. That is,

$$537 \quad (5.2) \quad \text{exec}(t_i) = \text{npic}^2(c_i) \quad \text{size}(d_j) = \text{npic}(c_j),$$

538 where $\text{npic}(c_i)$ denotes the number of particles in cell c_i .

539 We did not include matrices with less than $100 \times K$ matrix rows so that each
 540 processor will be assigned at least 100 rows on average. That is, we have matrices
 541 that have at least 6400 rows for 64-way partitioning and 12800 rows for 128-way
 542 partitioning. As a result of this selection criterion, the experiments are conducted for
 543 a total of 464 partitioning instances (117, 108, 94, 82, and 63 instances for 64-, 128-,
 544 256-, 512-, and 1024-way partitions, respectively).

545 **5.3.2. SpGEMM partitioning instances.** Consider the SpGEMM applica-
 546 tion of the form $C=AB$, where input matrices A and B are of sizes $q \times r$ and $r \times s$. In
 547 row-row-parallel SpGEMM, the atomic computational task t_i is the pre-multiplication
 548 of row i of matrix A with the whole matrix B . Then, the task set $\mathcal{T}=\{t_1, t_2, \dots, t_q\}$

549 contains q tasks $\{r_1^A B, r_2^A B, \dots, r_q^A B\}$, where r_i^A denotes row i of A matrix. The
 550 sparse-vector-matrix multiplication $r_i^A B$ requires row r_i^A as well as those B -matrix
 551 rows that correspond to the column indices (*cols*) of the nonzeros of row r_i^A . That is,

$$552 \quad (5.3) \quad \text{Data}(t_i) = \{a_{i,*}\} \cup \{b_{x,*} : x \in \text{cols}(a_{i,*})\}.$$

553 So, the sets of A - and B -matrix rows constitute the set of $q + r$ data elements. That
 554 is, $\mathcal{D} = \{d_1, d_2, \dots, d_{q+r}\} = \{r_1^A, r_2^A, \dots, r_q^A, r_1^B, r_2^B, \dots, r_r^B\}$. Hence, this row-row-
 555 parallel SpGEMM application falls within the framework in [section 2](#), since vector-
 556 matrix multiplications are independent but they share B -matrix rows. The computa-
 557 tional costs for atomic tasks and sizes for data elements are easily defined as follows:

$$558 \quad (5.4) \quad \text{exec}(t_i) = \sum_{x \in \text{cols}(a_{i,*})} \text{nnz}(b_{x,*}),$$

$$559 \quad (5.5) \quad \text{size}(d_j) = \text{nnz}(a_{j,*}) \text{ for } 1 \leq j \leq q,$$

$$560 \quad (5.6) \quad \text{size}(d_j) = \text{nnz}(b_{j-q,*}) \text{ for } q + 1 \leq j \leq q + r.$$

561 Here, we consider two types of SpGEMM instances: $C = AB$ and $C = AA$. For $C =$
 562 AB , we generate 69 instances from the SuiteSparse matrix collection [\[23\]](#) in a similar
 563 way to [\[3\]](#). Matrices *amazon0302* and *amazon0312* are used as A matrices which
 564 represent the similarity between items and B matrices are generated utilizing a Zipf
 565 distribution (with exponent set to 3.0) to determine the item preferences and a uniform
 566 distribution to determine the users that prefer a specific item [\[31\]](#). In this setting,
 567 $C = AB$ gives the candidate items to be recommended to each user. We generated
 568 66 instances by considering the setup phase of Algebraic Multigrid methods [\[8\]](#) which
 569 involves the Galerkin product of the form RAP that necessitates two consecutive
 570 SpGEMM operations. 11 of these instances are of the form $C = RA$, whereas the
 571 remaining 55 instances are of the form $C = AP$. The last instance in this category
 572 contains two different matrices, namely *thermomech_dK* and *thermomech_dM*, which
 573 are conformable for multiplication.

574 For the $C = AA$ type of instances, we selected 12 matrices from the SuiteSparse
 575 Matrix Collection [\[23\]](#). The number of rows/columns and number of nonzeros are
 576 in the range of 88K - 1.5M, and 2.5M - 30M, respectively. For the $C = AA$ type
 577 of instances, the fact that the B matrix is actually the A matrix can be exploited
 578 in order to reduce the memory footprint of the application. On the other hand, for
 579 the sake of computational efficiency, our parallel SpGEMM implementation does not
 580 exploit this fact. That is, we partition $C = AA$ instances as we partition $C = AB$
 581 instances.

582 The same partitioning granularity principle utilized for mesh instances is also
 583 used for SpGEMM instances. So, experimental results of the $C = AB$ instances are
 584 reported for a total of 259 partitioning instances (69, 63, 54, 45, and 28 instances for
 585 64-, 128-, 256-, 512-, and 1024-way partitions, respectively). Experimental results of
 586 the $C = AA$ instances are reported for a total of 59 partitioning instances (12, 12, 12,
 587 12, and 11 instances for 64-, 128-, 256-, 512-, and 1024-way partitions, respectively).

588 **5.3.3. Data size variation.** Here, we compare and discuss the irregularity of
 589 the datasets in terms of coefficient of variation (CV) values on the data sizes, where
 590 CV values are computed as standard deviation divided by mean. The purpose is to
 591 observe the relation between load balancing performance of the proposed algorithms
 592 and the irregularity of the datasets defined as the size variation of the data elements.

593 **Table 1** displays average CV values as well as the number of partitioning instances
 594 for each dataset and for each number of processors. Here higher CV values correspond
 595 to higher irregularity of the data sizes. As seen in the table, the $C = AA$ dataset has
 596 the largest average CV (0.90 on $K = 64$ processors), the $C = AB$ dataset has the
 597 smallest average CV (0.32 on $K = 64$ processors), and the mesh dataset has the in-
 598 between average CV (0.68 on $K = 64$ processors). So, in terms of data size variation,
 599 the $C = AA$ dataset is the most irregular dataset, whereas the $C = AB$ and the mesh
 600 datasets are respectively least and in-between irregular datasets.

TABLE 1
 Number of partitioning instances and average coefficient of variation (CV) values for each dataset

K	mesh dataset		$C = AB$ dataset		$C = AA$ dataset	
	# of ins.	CV	# of ins.	CV	# of ins.	CV
64	117	0.68	69	0.32	12	0.90
128	108	0.66	64	0.34	12	0.90
256	94	0.65	54	0.34	12	0.90
512	82	0.65	45	0.35	12	0.90
1024	63	0.60	28	0.36	11	0.81

601 5.4. Performance comparison.

602 **5.4.1. Performance metrics.** For each application, we evaluate the perfor-
 603 mance of the partitioning models in terms of maximum computational load CL_{max}
 604 and maximum data load DL_{max} handled by a processor (given in (2.4) and (2.5),
 605 respectively). For a simpler presentation, we give ratios of those metrics to their
 606 averages. That is,

$$607 \quad (5.7) \quad CL_{max}^r = \frac{CL_{max}}{CL_{avg}}, \text{ where } CL_{avg} = \frac{1}{K} CL_{tot} = \frac{1}{K} \sum_{t_i \in \mathcal{T}} exec(t_i),$$

$$608 \quad (5.8) \quad DL_{max}^r = \frac{DL_{max}}{DL_{avg}^*}, \text{ where } DL_{avg}^* = \frac{1}{K} DL_{tot} = \frac{1}{K} \sum_{d_j \in \mathcal{D}} size(d_j).$$

610 Here, CL_{avg} denotes the average computational load per processor under perfect load
 611 balance. DL_{avg}^* denotes average data load per processor without data replication
 612 under perfect load balance. So, DL_{avg}^* denotes the ideal average data load.

613 We also report the total data size (original total size together with replicated data
 614 size) incurred by the models as the ratio to original total data size as follows:

$$615 \quad (5.9) \quad DL_{rep}^r = \frac{\sum_{k=1}^K DL(p_k)}{DL_{tot}} = \frac{\sum_{k=1}^K \sum_{d_j \in WS(p_k)} size(d_j)}{DL_{tot}}.$$

616 This metric also defines a lower bound for the DL_{max}^r metric. That is, under perfect
 617 data load balance $DL_{max}^r = DL_{rep}^r$. It is clear that, in each of these three metrics, a
 618 smaller value refers to a better performance.

619 Recall that both proposed models \mathcal{H}_{DV} and \mathcal{H}_{IW} utilize the RB framework to in-
 620 crease their effectiveness. So, we also report the performance results for these models
 621 without utilizing the RB framework in order to show the relative performance improve-
 622 ment attained by the use of the RB framework. These experiments are performed by

623 directly K -way partitioning the hypergraph constructed at the very beginning. We
 624 use $_{RB}$ and $_{Dr}$ subscripts to refer to models which utilize RB and which do not utilize
 625 RB, respectively. That is, $\mathcal{H}_{DV_{RB}}$ and $\mathcal{H}_{IW_{RB}}$ refer to the models that utilize the RB
 626 framework, whereas $\mathcal{H}_{DV_{Dr}}$ and $\mathcal{H}_{IW_{Dr}}$ refer to the models that do utilize the RB
 627 framework.

628 **5.4.2. Average performance comparison.** Table 2 displays average perfor-
 629 mance comparison for the mesh partitioning instances. As seen in the table, both pro-
 630 posed models perform much better data load balancing than the baseline model with-
 631 out any performance degradation in computational load balance. The performance
 632 gap between the proposed models and the baseline model increases with increasing
 633 number of processors in favor of the proposed models. For example, performance
 634 improvement of $\mathcal{H}_{DV_{RB}}$ over \mathcal{H}_{Base} in the DL_{max}^r metric is 28%, 30%, 34%, 39%, and
 635 45% on respectively $K \in \{64, 128, 256, 512, 1024\}$ processors, .

TABLE 2
 Average performance comparison for the mesh partitioning instances

K	number of instances	\mathcal{H}_{Base}	$\mathcal{H}_{DV_{Dr}}$	$\mathcal{H}_{DV_{RB}}$	$\mathcal{H}_{IW_{Dr}}$	$\mathcal{H}_{IW_{RB}}$
		DL_{max}^r : maximum data load ratio				
64	117	1.83	1.33	1.31	1.29	1.25
128	108	2.02	1.44	1.41	1.40	1.34
256	94	2.28	1.56	1.51	1.50	1.42
512	82	2.59	1.69	1.58	1.61	1.49
1024	63	2.88	1.71	1.59	1.66	1.50
		DL_{rep}^r : total replication ratio				
64	117	1.11	1.18	1.16	1.18	1.16
128	108	1.15	1.24	1.21	1.24	1.22
256	94	1.18	1.29	1.26	1.29	1.27
512	82	1.21	1.34	1.29	1.35	1.31
1024	63	1.18	1.33	1.26	1.33	1.28
		CL_{max}^r : maximum computational load ratio				
64	117	1.04	1.03	1.04	1.03	1.04
128	108	1.04	1.03	1.04	1.03	1.04
256	94	1.05	1.03	1.05	1.04	1.05
512	82	1.06	1.04	1.05	1.05	1.06
1024	63	1.06	1.03	1.05	1.04	1.06

636 As seen in Table 2, the proposed models considerably increase the total amount
 637 of data replication compared to the baseline model. This is expected since two-
 638 constraint formulation limits the search space during the partitioning. For example, on
 639 $K=1024$ processors, \mathcal{H}_{Base} incurs only 18% replication, whereas $\mathcal{H}_{DV_{RB}}$ and $\mathcal{H}_{IW_{RB}}$
 640 incur 26% and 28% replication, respectively. On the other hand, this increase is not
 641 important since the proposed models significantly reduce the load of the maximally
 642 loaded processor via much better data load balancing. The gap between the DL_{max}^r
 643 and DL_{rep}^r metrics, which shows how close the model approaches to the lower bound,
 644 is smaller for both proposed models compared to the baseline model. For example,
 645 on $K=1024$ processors, \mathcal{H}_{Base} achieves average DL_{max}^r value of 144% (2.88 versus
 646 1.18) above the lower bound determined by the DL_{rep}^r value, whereas $\mathcal{H}_{DV_{RB}}$ and

647 $\mathcal{H}_{IW_{RB}}$ respectively achieve DL_{max}^r values of only 26% (1.59 versus 1.26) and 17%
 648 (1.50 versus 1.28) above the lower bounds.

649 As also seen in the Table 2, the use of the RB framework leads to considerable per-
 650 formance improvement in both proposed models and this performance improvement
 651 increases with increasing number of processors as expected. For example, performance
 652 improvement of $\mathcal{H}_{IW_{RB}}$ over $\mathcal{H}_{IW_{Dr}}$ in the DL_{max}^r metric is 3.1%, 4.3%, 5.3%, 7.5%,
 653 and 9.6% respectively on $K \in \{64, 128, 256, 512, 1024\}$ processors.

TABLE 3
 Average performance comparison for the $C=AB$ SpGEMM partitioning instances

K	number of instances	\mathcal{H}_{Base}	$\mathcal{H}_{DV_{Dr}}$	$\mathcal{H}_{DV_{RB}}$	$\mathcal{H}_{IW_{Dr}}$	$\mathcal{H}_{IW_{RB}}$
		DL_{max}^r : maximum data load ratio				
64	69	1.19	1.12	1.10	1.12	1.09
128	64	1.23	1.16	1.13	1.15	1.11
256	54	1.25	1.19	1.14	1.19	1.13
512	45	1.28	1.24	1.17	1.22	1.15
1024	28	1.30	1.28	1.19	1.26	1.17
DL_{rep}^r : total replication ratio						
64	69	1.05	1.07	1.06	1.07	1.06
128	64	1.06	1.09	1.07	1.09	1.07
256	54	1.07	1.11	1.08	1.11	1.09
512	45	1.09	1.13	1.10	1.13	1.10
1024	28	1.09	1.14	1.10	1.14	1.11
CL_{max}^r : maximum computational load ratio						
64	69	1.03	1.02	1.02	1.02	1.03
128	64	1.03	1.02	1.03	1.02	1.03
256	54	1.03	1.02	1.02	1.02	1.03
512	45	1.03	1.02	1.02	1.02	1.03
1024	28	1.03	1.02	1.03	1.02	1.03

654 Table 3 displays average performance comparison for the $C=AB$ type of SpGEMM
 655 partitioning instances. As seen in the table, DL_{max}^r values for this dataset are much
 656 less than those for mesh partitioning instances, which is because of the considerably
 657 less data size irregularity of $C=AB$ type of SpGEMM partitioning instances com-
 658 pared to that of mesh partitioning instances. For example, on $K=1024$ processors,
 659 the DL_{max}^r values attained by the different models vary between 1.30 and 1.17 on
 660 the $C=AB$ type of SpGEMM partitioning instances, whereas those values for the
 661 mesh partitioning instances vary between 2.88 and 1.50. As seen in the table, the
 662 proposed models perform about 10% better compared to the baseline model in the
 663 DL_{max}^r metric. Similar to the mesh instances, use of the RB framework increases the
 664 performance of the proposed model, as expected. For example, on 1024 processors,
 665 $\mathcal{H}_{DV_{RB}}$ performs 7.0% better than the $\mathcal{H}_{DV_{Dr}}$ model, whereas $\mathcal{H}_{IW_{RB}}$ performs 7.1%
 666 better than the $\mathcal{H}_{IW_{Dr}}$ model. On the other hand, in the DL_{rep}^r metric the proposed
 667 models do not increase the total replication considerably, in contrast to the mesh par-
 668 titioning instances. Similar to the mesh partitioning instances, CL_{max}^r metric remains
 669 almost the same for all instances. Here, obtained improvement rates are significantly
 670 less than compared to the mesh partitioning instances. This can be attributed to
 671 regularity of the $C=AB$ type of SpGEMM instances.

TABLE 4
Average performance comparison for the $C=AA$ SpGEMM partitioning instances

K	number of instances	\mathcal{H}_{Base}	$\mathcal{H}_{DV_{Dr}}$	$\mathcal{H}_{DV_{RB}}$	$\mathcal{H}_{IW_{Dr}}$	$\mathcal{H}_{IW_{RB}}$
		DL_{max}^r : maximum data load ratio				
64	12	3.81	3.06	2.61	3.09	2.44
128	12	4.85	4.06	3.30	4.20	3.14
256	12	6.50	5.59	4.72	6.02	4.39
512	12	9.58	8.23	7.11	8.41	6.53
1024	11	12.51	10.88	9.07	11.24	8.91
		DL_{rep}^r : total replication ratio				
64	12	1.81	2.12	1.99	2.10	2.02
128	12	2.09	2.56	2.36	2.54	2.42
256	12	2.50	3.21	2.90	3.17	2.98
512	12	3.10	4.15	3.65	4.08	3.81
1024	11	3.86	5.20	4.59	5.12	4.83
		CL_{max}^r : maximum computational load ratio				
64	12	1.03	1.03	1.03	1.03	1.04
128	12	1.03	1.03	1.05	1.03	1.06
256	12	1.04	1.03	1.06	1.04	1.06
512	12	1.05	1.03	1.06	1.05	1.09
1024	11	1.06	1.03	1.05	1.04	1.13

672 Table 4 displays average performance comparison for the $C=AA$ type of SpGEMM
673 partitioning instances. Much higher irregularity in data size distribution of this
674 dataset incurs much higher DL_{max}^r values as seen in Tables 2 to 4. For example,
675 on $K=1024$, \mathcal{H}_{Base} obtains DL_{max}^r values of 12.51, 1.30, and 2.88 respectively for
676 $C=AA$, $C=AB$, and mesh partitioning instances. Such partitioning instances with
677 high data size variation incur hard partitioning instances and justify the importance
678 of the target optimization problem.

679 As seen in Table 4, both proposed models perform much better data load balanc-
680 ing than the baseline model with almost no performance degradation in computational
681 load balance. The computational load balance becomes considerably worse for $\mathcal{H}_{IW_{RB}}$
682 model only on 512 and 1024 processors. The proposed models $\mathcal{H}_{IW_{RB}}$ and $\mathcal{H}_{DV_{RB}}$
683 perform 27% and 29% better than \mathcal{H}_{Base} model on 1024 processors.

684 As seen in Table 4, the total amount of data replication increases considerably
685 compared to the baseline model. For example, on $K=1024$ processors, \mathcal{H}_{Base} in-
686 curs 286% replication, whereas $\mathcal{H}_{DV_{RB}}$ and $\mathcal{H}_{IW_{RB}}$ incur 359% and 383% replication,
687 respectively. The much higher difference in the DL_{max}^r and DL_{rep}^r values for this
688 dataset (for example, 9.07 versus 4.59 for $\mathcal{H}_{DV_{RB}}$ on 1024 processors) compared to
689 the other two datasets is because of the much higher data size variation in this dataset.

690 As also seen in Table 4, \mathcal{H}_{IW} benefits more from the RB framework compared
691 to \mathcal{H}_{DV} . For example, on $K=1024$ processors, the RB framework increases the
692 performance of \mathcal{H}_{DV} and \mathcal{H}_{IW} by 17% and 21%, respectively. This can be also
693 observed for mesh and $C=AB$ instances. For example, on $K=1024$ processors,
694 for the mesh dataset, RB framework increases the performance of \mathcal{H}_{DV} and \mathcal{H}_{IW} by
695 7.0% and 9.6%, respectively. Higher sensitivity of the use of the RB framework on the
696 \mathcal{H}_{DV} model on the $C=AA$ dataset makes $\mathcal{H}_{DV_{Dr}}$ model perform better than $\mathcal{H}_{IW_{Dr}}$,

697 although $\mathcal{H}_{IW_{RB}}$ performs better than $\mathcal{H}_{DV_{RB}}$ on average for each K value. This is
 698 because, as we discussed earlier (in [subsection 4.3](#)), the \mathcal{H}_{DV} model tries to encode
 699 data loads of processors after the RB step depending on the resulting bipartition,
 700 whereas the \mathcal{H}_{IW} model estimates data loads of processors before the RB step and
 701 then corrects its estimation depending on the resulting bipartition for the following
 702 RB level.

TABLE 5
 Number of instances for which each model attains the best performance on DL_{max}^r

K	mesh dataset			$C = AB$ dataset			$C = AA$ dataset		
	\mathcal{H}_{Base}	\mathcal{H}_{DV}	\mathcal{H}_{IW}	\mathcal{H}_{Base}	\mathcal{H}_{DV}	\mathcal{H}_{IW}	\mathcal{H}_{Base}	\mathcal{H}_{DV}	\mathcal{H}_{IW}
64	2	40	75	5	22	42	1	0	11
128	6	31	71	5	12	47	1	3	8
256	6	22	66	4	11	39	1	3	8
512	6	25	51	5	11	29	1	1	10
1024	3	18	42	5	5	18	0	2	9

703 The relative performance comparison of the two proposed models \mathcal{H}_{DV} and \mathcal{H}_{IW}
 704 is as follows: As seen in [Tables 2 to 4](#), \mathcal{H}_{IW} display slightly better average performance
 705 than \mathcal{H}_{DV} for each dataset and for each K value, except for the $C = AA$ dataset where
 706 $\mathcal{H}_{DV_{Dr}}$ performs slightly better than $\mathcal{H}_{IW_{Dr}}$ for each K value on average. [Table 5](#) is
 707 presented to show the number of partitioning instances for which each model attains
 708 the best performance in the DL_{max}^r metric. As seen in the table, although \mathcal{H}_{IW}
 709 attains the highest number of best DL_{max}^r values, \mathcal{H}_{DV} attains considerably high
 710 number of best DL_{max}^r values especially on the mesh and the $C = AB$ datasets. So
 711 both proposed models should be considered for attaining simultaneous computational
 712 and data load balancing depending on the nature of the application and the dataset.

713 **5.4.3. Performance profile comparison.** In [Figures 7 and 8](#), we provide per-
 714 formance profiles for each of the five models in terms of the DL_{max}^r and DL_{rep}^r metrics
 715 on $K \in \{64, 128, 256, 512, 1024\}$ processors. In each figure, each column corresponds
 716 to a different dataset, whereas each row corresponds to a different K value. We do
 717 not include performance profiles for the CL_{max}^r metric, since each of the five models
 718 perform similarly as discussed earlier. That is, performance profiles for the CL_{max}^r
 719 metric are almost on top of each other.

720 Performance profiles [28] are widely used in comparing multiple models over a
 721 large collection of test cases. In a different way from the average performance compar-
 722 ison given above, in the performance profiles we compare five models according to
 723 the best performing model for each partitioning instance on one of the metrics and
 724 measure in what fraction of the test cases a model performs within a factor of the
 725 best observed performance. A point (x, y) in a curve means that the respective model
 726 is within an x factor of the best result in a fraction y of the dataset. For example,
 727 consider the point $(x = 1.05, y = 0.70)$ on the performance curve for $\mathcal{H}_{DV_{RB}}$ for the
 728 mesh dataset on $K = 1024$ processors for the DL_{max}^r metric. This point means that
 729 for 70% of the partitioning instances $\mathcal{H}_{DV_{RB}}$ attains DL_{max}^r values at most 5 percent
 730 larger than the best DL_{max}^r values achieved. Thus, a model that is closer to the
 731 top-left corner is better.

732 As seen in [Figures 7 and 8](#), the performance profiles corroborate the above discus-
 733 sion about the relative performance of the models. For example, comparison of profile
 734 curves for $\mathcal{H}_{DV_{RB}}$ and $\mathcal{H}_{DV_{Dr}}$, as well as for $\mathcal{H}_{IW_{RB}}$ and $\mathcal{H}_{IW_{Dr}}$ on the DL_{max}^r metric

735 show that the performance improvement obtained from the use of the RB framework
 736 increases with increasing number of processors for the mesh and the $C = AB$ datasets,
 737 whereas this improvement remains almost the same for the $C = AA$ dataset. Further-
 738 more, comparison of profile curves for $\mathcal{H}_{DV_{Dr}}$ and $\mathcal{H}_{IW_{Dr}}$ on DL_{max}^r metric shows
 739 that the performance obtained by the $\mathcal{H}_{DV_{Dr}}$ becomes superior than the performance
 740 obtained by $\mathcal{H}_{IW_{Dr}}$ for about more than 20% of the $C = AA$ dataset within a factor
 741 of approximately 1.20.

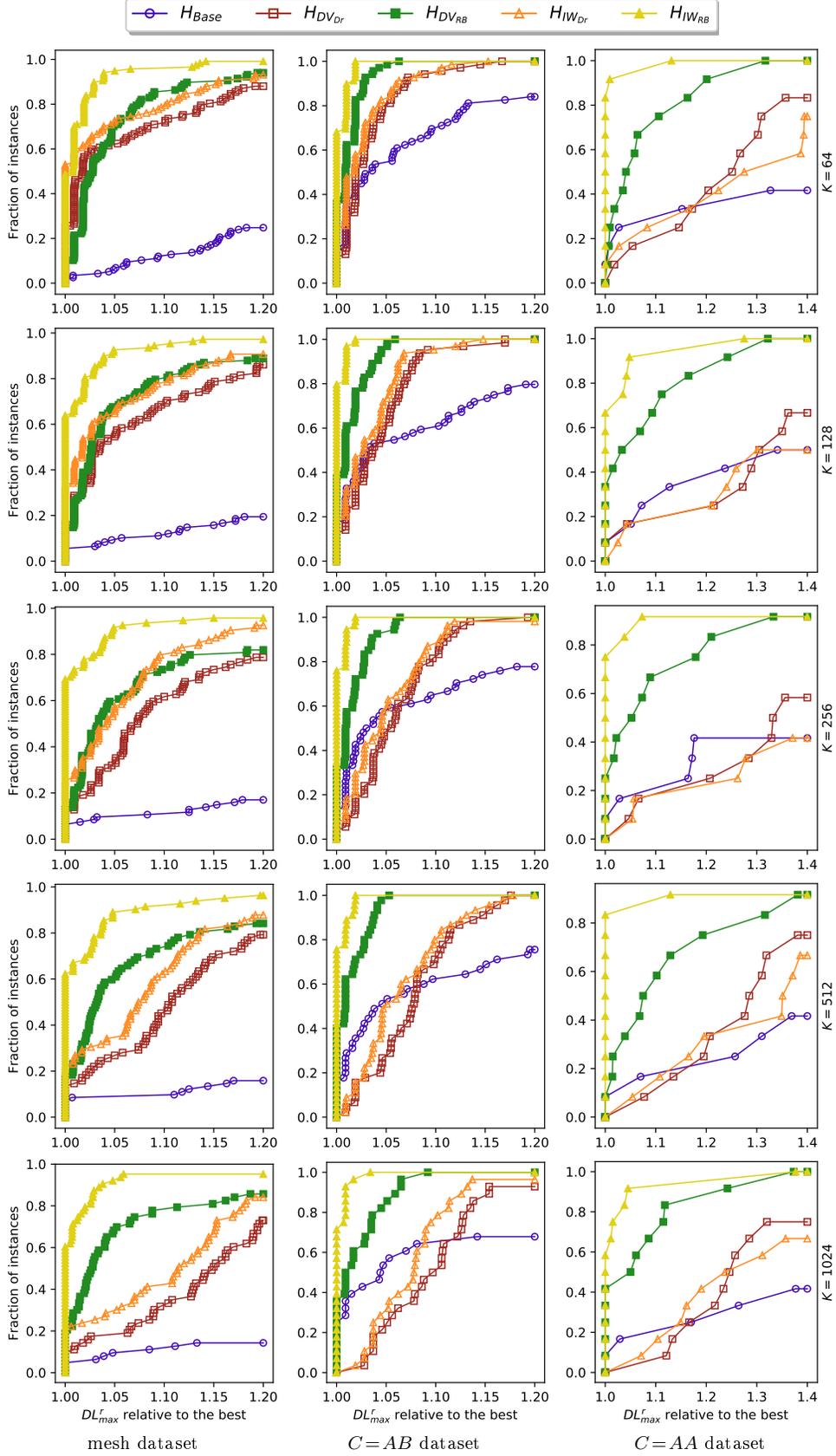
742 As seen in Figure 7, on $K = 1024$ processors, $\mathcal{H}_{IW_{RB}}$ achieves the best solutions for
 743 DL_{max}^r metric for the 60% of all datasets, whereas $\mathcal{H}_{DV_{RB}}$ attains solutions within
 744 factors of 1.04, 1.02, and 1.08 for the mesh, $C = AB$ and $C = AA$ partitioning in-
 745 stances, respectively. As seen in Figure 8, the \mathcal{H}_{Base} performs the best on the DL_{rep}^r
 746 metric, whereas $\mathcal{H}_{IW_{RB}}$ and $\mathcal{H}_{DV_{RB}}$ follows. Furthermore, this difference increases
 747 with increasing irregularity of the dataset.

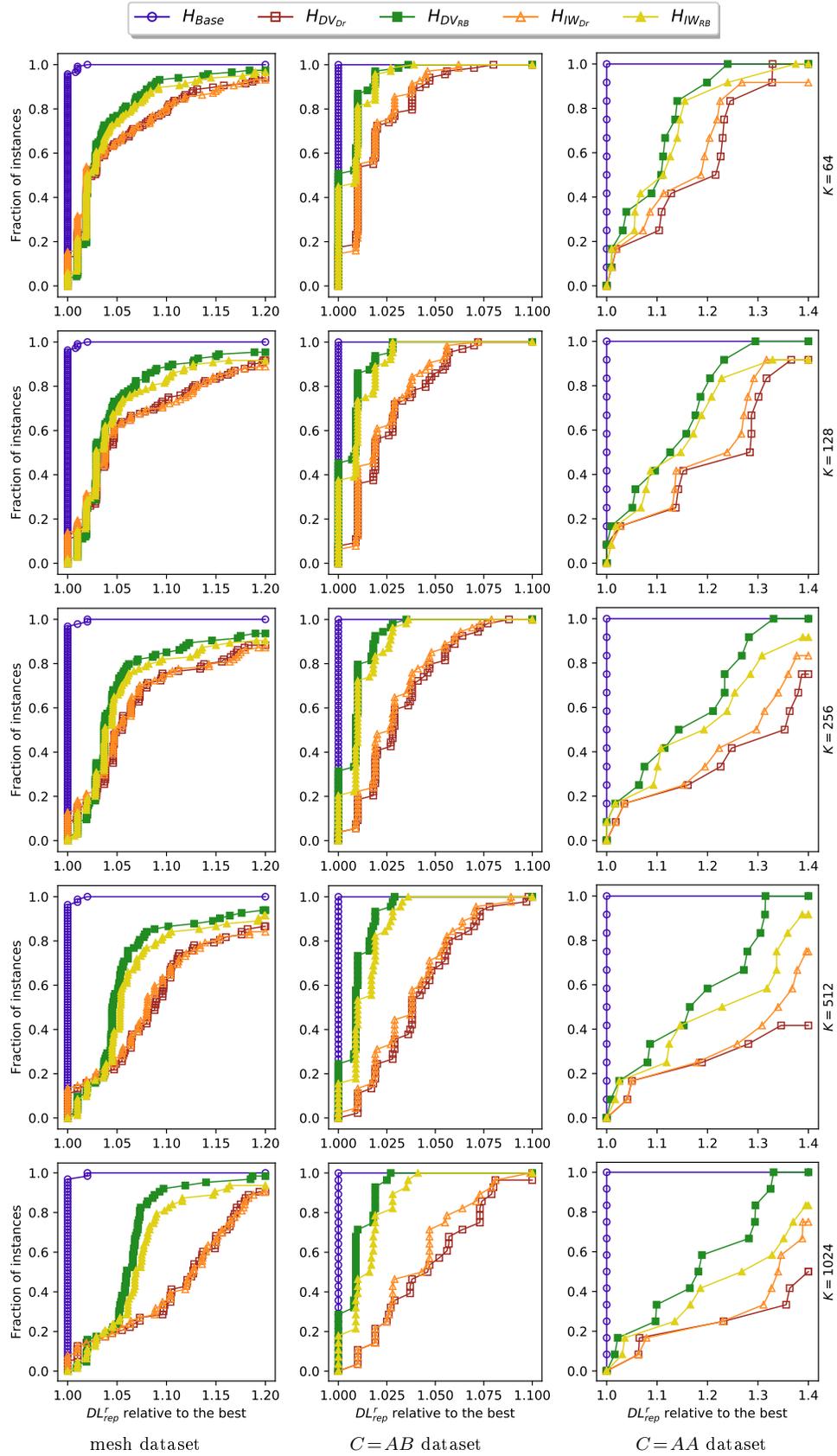
748 **6. Related work.** Chevalier et al. [22] consider memory constraints during the
 749 task partitioning with the objective of balancing computational load of processors.
 750 They propose a multilevel bipartite graph partitioning algorithm as follows: In the
 751 coarsening phase, computational vertices match with other computational vertices
 752 according to their data element share. Data elements are matched with other data
 753 elements similar to the identical net detection in HP. They perform initial partition-
 754 ing by greedily assigning sorted tasks to parts/processors considering the memory
 755 capacity. In the refinement phase, tasks are moved between parts/processors with the
 756 objective of decreasing the computational load. Note that this work does not balances
 757 the data load rather abides by the memory constraint.

758 Angel et al. [4, 5] also consider memory constraints while minimizing the com-
 759 putational load of the maximally loaded processor. In [4], a dynamic programming
 760 based Fixed-Parameter Tractable algorithm with respect to the path-width of the
 761 neighborhood graph is utilized. Note that this is an approximation since the problem
 762 is NP-hard. In [5], they focus on the case when the neighborhood graph has bounded
 763 tree-width. In that way tree decomposition of the graph and its traversal in a specific
 764 way which may be useful on its own allows to find a solution within a factor of epsilon.

765 Tzovas et al. [40] target heterogeneous distributed systems for which sparse ma-
 766 trices or graphs will be distributed. Also in this work, memory capacity is a constraint
 767 where the objective is maximizing utilization of each processing unit with minimum
 768 communication cost. For this purpose, they propose a two phase algorithm: In the
 769 first phase, they greedily compute block size for each processing unit. In the second
 770 phase, they feed these sizes to variety a of existing partitioning tools.

771 **7. Conclusion.** In high performance computing applications, computational
 772 load balancing is well studied in the literature, however data load balancing has be-
 773 come a concern recently. In that sense, to our knowledge, this is a pioneer work
 774 in which we simultaneously balance processors' computational and data loads in
 775 distributed-memory setting. We proposed two different hypergraph partitioning based
 776 models both utilizing a two-constraint formulation for load balancing. The first con-
 777 straint encodes balancing the computational load, whereas the second constraint en-
 778 codes balancing the data load. The partitioning objective encodes minimizing the
 779 total amount of data replication. We utilized the well known recursive bipartition-
 780 ing framework for increasing the accuracy of the both proposed models. Instead of
 781 developing a new partitioner from scratch both proposed models can easily be im-
 782 plemented by invoking any HP tool that supports multi constraint partitioning as
 783 a two-way partitioner at each RB step. Extensive experiments show that both pro-

FIG. 7. Performance profiles for the DL_{max}^r metric on $K \in \{64, 128, 256, 512, 1024\}$

FIG. 8. Performance profiles for the DL_{rep}^r metric on $K \in \{64, 128, 256, 512, 1024\}$

784 posed models perform significantly better than a baseline model. We achieve up to
785 49 percent better performance on 1024 processors.

786

REFERENCES

- 787 [1] S. ACER, O. SELVITOPİ, AND C. AYKANAT, *Improving performance of sparse matrix dense*
788 *matrix multiplication on large-scale parallel systems*, *Parallel Computing*, 59 (2016), pp. 71
789 – 96. *Theory and Practice of Irregular Applications*.
- 790 [2] S. ACER, O. SELVITOPİ, AND C. AYKANAT, *Optimizing nonzero-based sparse matrix parti-*
791 *tioning models via reducing latency*, *Journal of Parallel and Distributed Computing*, 122
792 (2018), pp. 145 – 158.
- 793 [3] K. AKBUDAK, O. SELVITOPİ, AND C. AYKANAT, *Partitioning models for scaling parallel sparse*
794 *matrix-matrix multiplication*, *ACM Trans. Parallel Comput.*, 4 (2018).
- 795 [4] E. ANGEL, C. CHEVALIER, F. LEDOUX, S. MORAIS, AND D. REGNAULT, *FPT approximation*
796 *algorithm for scheduling with memory constraints*, in *Euro-Par 2016: Parallel Processing*,
797 Cham, 2016, Springer International Publishing, pp. 196–208.
- 798 [5] E. ANGEL, S. MORAIS, AND D. REGNAULT, *A bi-criteria FPTAS for scheduling with memory*
799 *constraints on graph with bounded tree-width*, arXiv preprint arXiv:2202.08704, (2022).
- 800 [6] C. AYKANAT, B. B. CAMBAZOGLU, F. FINDIK, AND T. KURC, *Adaptive decomposition and*
801 *remapping algorithms for object-space-parallel direct volume rendering of unstructured*
802 *grids*, *Journal of Parallel and Distributed Computing*, 67 (2007), pp. 77–99.
- 803 [7] R. BARAT, *Load Balancing of Multi-physics Simulation by Multi-criteria Graph Partition-*
804 *ing*, PhD thesis, Thèse de doctorat dirigée par Pellegrini, François et Chevalier, Cédric
805 Informatique Bordeaux 2017, 2017. 2017BORD0961.
- 806 [8] N. BELL, S. DALTON, AND L. N. OLSON, *Exposing fine-grained parallelism in algebraic multi-*
807 *grid methods*, *SIAM Journal on Scientific Computing*, 34 (2012), pp. C123–C152.
- 808 [9] R. H. BISSELING, *Parallel Scientific Computation: A Structured Approach using BSP (Second*
809 *edition)*, Oxford University Press, USA, 2020.
- 810 [10] R. H. BISSELING AND I. FLESCH, *Mondriaan sparse matrix partitioning for attacking cryp-*
811 *tosystems by a parallel block Lanczos algorithm – a case study*, *Parallel Computing*, 32
812 (2006), pp. 551 – 567. *Algorithmic Skeletons*.
- 813 [11] R. H. BISSELING AND W. MEESEN, *Communication balancing in parallel sparse matrix-vector*
814 *multiplication*, *Electronic Transactions on Numerical Analysis*, 21 (2005), pp. 47–65.
- 815 [12] E. G. BOMAN AND M. M. WOLF, *A nested dissection partitioning method for parallel sparse*
816 *matrix-vector multiplication*, in *2013 IEEE High Performance Extreme Computing Con-*
817 *ference (HPEC)*, 2013, pp. 1–6.
- 818 [13] B. B. CAMBAZOGLU AND C. AYKANAT, *Hypergraph-partitioning-based remapping models for*
819 *image-space-parallel direct volume rendering of unstructured grids*, *IEEE Transactions on*
820 *Parallel and Distributed Systems*, 18 (2007), pp. 3–16.
- 821 [14] Ü. V. ÇATALYÜREK AND C. AYKANAT, *Hypergraph-partitioning-based decomposition for par-*
822 *allel sparse-matrix vector multiplication*, *IEEE Transactions on Parallel and Distributed*
823 *Systems*, 10 (1999), pp. 673–693.
- 824 [15] Ü. V. ÇATALYÜREK AND C. AYKANAT, *A hypergraph-partitioning approach for coarse-grain*
825 *decomposition*, in *Proceedings of the 2001 ACM/IEEE Conference on Supercomputing*, SC
826 '01, New York, NY, USA, 2001, Association for Computing Machinery, p. 28.
- 827 [16] Ü. V. ÇATALYÜREK AND C. AYKANAT, *PaToH (partitioning tool for hypergraphs)*, in *Encyclo-*
828 *pedia of Parallel Computing*, Springer, 2011, pp. 1479–1487.
- 829 [17] Ü. V. ÇATALYÜREK, C. AYKANAT, AND B. UÇAR, *On two-dimensional sparse matrix parti-*
830 *tioning: Models, methods, and a recipe*, *SIAM Journal on Scientific Computing*, 32 (2010),
831 pp. 656–683.
- 832 [18] Ü. V. ÇATALYÜREK, E. G. BOMAN, K. D. DEVINE, D. BOZDAG, R. HEAPHY, AND L. A.
833 RIESEN, *Hypergraph-based dynamic load balancing for adaptive scientific computations*, in
834 *2007 IEEE International Parallel and Distributed Processing Symposium*, 2007, pp. 1–11.
- 835 [19] Ü. V. ÇATALYÜREK, E. G. BOMAN, K. D. DEVINE, D. BOZDAÇ, R. T. HEAPHY, AND L. A.
836 RIESEN, *A repartitioning hypergraph model for dynamic load balancing*, *Journal of Parallel*
837 *and Distributed Computing*, 69 (2009), pp. 711–724.
- 838 [20] Ü. V. ÇATALYÜREK, M. DEVECİ, K. KAYA, AND B. UÇAR, *UMPa: A multi-objective, multi-*
839 *level partitioner for communication minimization*, *Graph Partitioning and Graph Cluster-*
840 *ing*, 588 (2013), p. 53.
- 841 [21] A. CEVAHIR, C. AYKANAT, A. TURK, AND B. B. CAMBAZOGLU, *Site-based partitioning and*
842 *repartitioning techniques for parallel pagerank computation*, *IEEE Transactions on Parallel*

- 843 and Distributed Systems, 22 (2011), pp. 786–802.
- 844 [22] C. CHEVALIER, F. LEDOUX, AND S. MORAIS, *A Multilevel Mesh Partitioning Algorithm Driven*
845 *by Memory Constraints*, pp. 85–95.
- 846 [23] T. A. DAVIS AND Y. HU, *The University of Florida sparse matrix collection*, ACM Trans.
847 Math. Softw., 38 (2011).
- 848 [24] M. DEVECI, K. KAYA, B. UÇAR, AND Ü. V. ÇATALYÜREK, *Hypergraph partitioning for mul-*
849 *tiple communication cost metrics: Model and methods*, Journal of Parallel and Distributed
850 Computing, 77 (2015), pp. 69 – 83.
- 851 [25] K. DEVINE, B. HENDRICKSON, E. BOMAN, M. ST. JOHN, AND C. VAUGHAN, *Design of*
852 *dynamic load-balancing tools for parallel applications*, in Proceedings of the 14th Interna-
853 tional Conference on Supercomputing, ICS '00, New York, NY, USA, 2000, Association for
854 Computing Machinery, p. 110–118.
- 855 [26] K. D. DEVINE, E. G. BOMAN, R. T. HEAPHY, B. A. HENDRICKSON, J. D. TERESCO,
856 J. FAIK, J. E. FLAHERTY, AND L. G. GERVASIO, *New challenges in dynamic load balanc-*
857 *ing*, Applied Numerical Mathematics, 52 (2005), pp. 133–152. ADAPT '03: Conference on
858 Adaptive Methods for Partial Differential Equations and Large-Scale Computation.
- 859 [27] K. D. DEVINE, E. G. BOMAN, AND G. KARYPIS, *Partitioning and load balancing for emerging*
860 *parallel applications and architectures*, in Parallel Processing for Scientific Computing,
861 Society for Industrial and Applied Mathematics, 2006, pp. 99–126.
- 862 [28] E. D. DOLAN AND J. J. MORÉ, *Benchmarking optimization software with performance profiles*,
863 Mathematical programming, 91 (2002), pp. 201–213.
- 864 [29] J. FLAHERTY, R. LOY, P. SCULLY, M. SHEPHARD, B. SZYMANSKI, J. TERESCO, AND
865 L. ZIANTZ, *Load balancing and communication optimization for parallel adaptive finite*
866 *element methods*, in Proceedings 17th International Conference of the Chilean Computer
867 Science Society, 1997, pp. 246–255.
- 868 [30] S. JUN LIANG, J. CHENG, AND J. WEI ZHANG, *Research on data load balancing technology*
869 *of massive storage system for wearable devices*, Digital Communications and Networks,
870 (2020).
- 871 [31] G. LINDEN, B. SMITH, AND J. YORK, *Amazon.com recommendations: item-to-item collabo-*
872 *rative filtering*, IEEE Internet Computing, 7 (2003), pp. 76–80.
- 873 [32] K. LIU, G. XU, ET AL., *An improved hadoop data load balancing algorithm*, Journal of Net-
874 works, 8 (2013), p. 2816.
- 875 [33] H. MEYERHENKE, *Dynamic load balancing for parallel numerical simulations based on repar-*
876 *tioning with disturbed diffusion*, in 2009 15th International Conference on Parallel and
877 Distributed Systems, 2009, pp. 150–157.
- 878 [34] S. MORAIS, *Etude et obtention d'heuristiques et d'algorithmes exacts et approchés pour un*
879 *problème de partitionnement de maillage sous contraintes mémoire*, theses, Université
880 Paris Saclay, Nov. 2016.
- 881 [35] N. PATEL AND S. CHAUHAN, *A survey on load balancing and scheduling in cloud computing*,
882 Int. Journal for Innovative Research in Science and Technology, 1 (2015), pp. 185–189.
- 883 [36] D. M. PELT AND R. H. BISSELING, *A medium-grain method for fast 2D bipartitioning of sparse*
884 *matrices*, in 2014 IEEE 28th International Parallel and Distributed Processing Symposium,
885 May 2014, pp. 529–539.
- 886 [37] A. PINAR AND B. HENDRICKSON, *Improving load balance with flexibly assignable tasks*, IEEE
887 Transactions on Parallel and Distributed Systems, 16 (2005), pp. 956–965.
- 888 [38] O. SELVITOPİ, S. ACER, AND C. AYKANAT, *A recursive hypergraph bipartitioning framework*
889 *for reducing bandwidth and latency costs simultaneously*, IEEE Transactions on Parallel
890 and Distributed Systems, 28 (2017), pp. 345–358.
- 891 [39] O. SELVITOPİ AND C. AYKANAT, *Reducing latency cost in 2D sparse matrix partitioning models*,
892 Parallel Computing, 57 (2016), pp. 1 – 24.
- 893 [40] C. TZOVAS, M. PREDARI, AND H. MEYERHENKE, *Distributing sparse matrix/graph applica-*
894 *tions in heterogeneous clusters - an experimental study*, in 2020 IEEE 27th International
895 Conference on High Performance Computing, Data, and Analytics (HiPC), 2020, pp. 72–81.
- 896 [41] B. UÇAR AND C. AYKANAT, *Minimizing communication cost in fine-grain partitioning of*
897 *sparse matrices*, in International Symposium on Computer and Information Sciences,
898 Springer, 2003, pp. 926–933.
- 899 [42] B. UÇAR AND C. AYKANAT, *Encapsulating multiple communication-cost metrics in partition-*
900 *ing sparse rectangular matrices for parallel matrix-vector multiplies*, SIAM Journal on
901 Scientific Computing, 25 (2004), pp. 1837–1859.
- 902 [43] B. VASTENHOUW AND R. H. BISSELING, *A two-dimensional data distribution method for par-*
903 *allel sparse matrix-vector multiplication*, SIAM Review, 47 (2005), pp. 67–95.