PARTITIONING AND REORDERING FOR SPIKE-BASED DISTRIBUTED-MEMORY PARALLEL GAUSS-SEIDEL *

TUGBA TORUN[†], F. SUKRU TORUN[‡], MURAT MANGUOGLU[§], AND CEVDET 3 AYKANAT[†]

Abstract. Gauss-Seidel (GS) is a widely-used iterative method for solving sparse linear system 5 6 of equations and also known to be effective as a smoother in algebraic multigrid methods. Paralleliza-7 tion of GS is a challenging task since solving the sparse lower triangular system in GS constitutes a sequential bottleneck at each iteration. We propose a distributed-memory parallel GS (dmpGS) 8 9 by implementing a parallel sparse triangular solver (stSpike) based on the Spike algorithm. stSpike decouples the global triangular system into smaller systems that can be solved concurrently and 10 11 requires the solution of a much smaller reduced sparse lower triangular system which constitutes a sequential bottleneck. In order to alleviate this bottleneck and to reduce the communication over-12 13head of dmpGS, we propose a partitioning and reordering model consisting of two phases. The first 14phase is a novel hypergraph partitioning model whose partitioning objective simultaneously encodes 15 minimizing the reduced system size and the communication volume. The second phase is an in-block row reordering method for decreasing the nonzero count of the reduced system. Extensive experiments on a dataset consisting of 359 sparse linear systems verify the effectiveness of the proposed 17 partitioning and reordering model in terms of reducing the communication and the sequential com-18 19putational overheads. Parallel experiments on 12 large systems using up to 320 cores demonstrate that the proposed model significantly improves the scalability of dmpGS. 20

21 Key words. parallel Gauss-Seidel, distributed-memory, Spike algorithm, parallel sparse trian-22 gular solve, linear system solution, hypergraph partitioning, sparse matrix reordering.

AMS subject classifications. 68W10, 05C65, 65F50, 65F10 23

24 **1.** Introduction. A wide range of applications in science and engineering require the solution of a sparse linear system of equations 25

$$26 \quad (1.1) \qquad \qquad Ax = f,$$

1 2

4

where $A \in \mathbb{R}^{m \times m}$ is a general large sparse invertible matrix; and x and $f \in \mathbb{R}^m$ are 27the unknown and right hand side vectors, respectively. Depending on the numerical 28 and structural properties of the coefficient matrix, various solvers have been proposed. 29 Direct solvers require a sequence of operations: reordering and partitioning, sym-30 bolic factorization, numerical factorization, and finally obtaining the solution, typ-32 ically via forward and backward sweeps. The reordering and partitioning schemes are used both to reduce the amount of fill-in and to enhance the parallel scalability. Symbolic factorization is used to determine the sparsity structure of the factors, and 34 finally the numerical factorization (such as sparse LU [23], QR [34], SVD [13] and WZ [17]) is computed. Direct solvers are robust and, in general, are known to be very 36 scalable during the factorization phase [5, 43], but not so much during the triangular solution phase [45]. 38

^{*}Submitted to the editors April 11, 2021.

Funding: Computing resources used in this work were provided by the National Center for High Performance Computing of Turkey (UHeM) under grant number 4007772020. The third author was supported by the BAGEP Award of the Science Academy.

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

[‡]Department of Computer Engineering, Ankara Yildirim Beyazit University, 06010, Ankara, Turkey (ftorun@ybu.edu.tr).

[§]Department of Computer Engineering, Middle East Technical University, 06800 Ankara, Turkey (manguoglu@ceng.metu.edu.tr).

Iterative solvers, on the other hand, are known to be more scalable but not as 39 40 robust as direct solvers. Nevertheless, they are still preferred for large sparse systems due to their lower memory requirements. Starting with an initial guess for the solution 41 vector, these methods improve the solution at each iteration. There are two main types 42 of iterative solvers: stationary and non-stationary methods. Stationary methods have 43 the general form $x^{(k+1)} = \phi(x^{(k)})$ where $x^{(k)}$ is the solution vector at the k^{th} iteration 44 and $\phi(\cdot)$ is a function which does not change during the iterations. Some examples 45are Jacobi, Gauss-Seidel, Successive Over Relaxation (SOR) and Symmetric SOR 46 (SSOR) [34, 59]. Non-stationary methods have the form $x^{(k+1)} = \phi^{(k)}(x^{(k)})$ in which 47 the function $\phi^{(k)}(\cdot)$ changes at each iteration. Some examples are projection methods, 48 Krylov subspace methods and Chebyshev iterations [9, 35, 59]. 49

50 In practice, linear systems are preconditioned to reduce the required number of iterations of the iterative solvers and to improve their robustness. There could be a variety of choices of preconditioners, some are problem specific and others are more 52general. General classical preconditioners include, incomplete factorization based pre-53 conditoners (such as incomplete LU [58, 59]), sparse approximate inverse [11], alge-54braic multigrid (AMG) [51, 57], and others. We refer the reader to [10] for a detailed survey of preconditioners. Among these preconditioners, AMG has been widely used 56 recently in many applications [12, 30, 53] which is a generalization of Geometric Multigrid (GMG) [70]. GMG requires some knowledge of the physical problem and/or its 58 geometry, while there is no such requirement for AMG. AMG can be also used as a direct solver [36, 71]. Furthermore, AMG typically uses another iterative method as 61 a "smoother" which is required to reduce the error at each level and the smoother itself can also be preconditioned. More recently a preferred smoother for AMG is 62 Gauss-Seidel [3, 16, 67], as in BoomerAMG [36] and Trilinos-ML [32]. 63

Gauss-Seidel (GS) is a well-known stationary iterative method which solves the linear system (1.1) by splitting the coefficient matrix into its lower and strictly upper triangular parts, A = L+U. Then the solution is obtained iteratively by

67
$$x^{(k+1)} = L^{-1}(f - Ux^{(k)}).$$

At each iteration of GS, both a lower triangular system is required to be solved and an upper triangular SpMV (sparse matrix-vector multiplication) is performed. GS is guaranteed to converge if A is strictly or irreducibly diagonal dominant [7] or symmetric positive definite [34]. It is known to be effective and preferred as a smoother for a wide variety of problems [3, 72]. However, a true distributed-memory parallelization of GS is considered to be a challenging task [3].

In the literature, parallel GS implementations are proposed either to solve the 74 original problem (1.1) [6, 42, 62] or to use it as a smoother in multigrid schemes 76 [38, 64, 73]. A commonly-used method to parallelize GS by finding independent sub-tasks is the red-black coloring strategy [2, 31, 41], which has been extended to 77 multi-coloring [33, 52, 4] to attain more parallelism for complicated regular problems. 78 However, multi-colored GS is not feasible for some cases such as unstructured finite 79 element applications since the number of colors becomes too large [42]. Another 80 81 approach is to use a processor-localized GS in which each processor performs GS as a subdomain solver, but its convergence rate is low and may diverge for a large number 82 83 of processors [3].

The main difficulty in parallelizing GS inherits from the sequential nature of triangular solve included in GS [72]. Along with its importance in several applications, solving triangular systems often constitutes a sequential bottleneck because of the dependencies between unknowns in forward or backward substitutions. In [60], a

parallel banded triangular solver is proposed. This algorithm is extended for solving 88 banded linear systems [21, 28] and further improved by implementing various alterna-89 tives in each step of the factorization including the solution of the reduced system in 90 [55, 56, 63]. At this point, the algorithm is called Spike algorithm. For sparse linear 91 systems, Spike is also proposed as a solver for a banded preconditioner that is sparse within the band [49, 61], and it is generalized for sparse linear systems [15, 47, 48]. 93 In [69], a Spike-based parallel solver for general tridiagonal systems is implemented 94 for GPU architectures. A recent study [22] proposes a multi-threaded parallel solver 95 for sparse triangular systems by extending the Spike algorithm [60]. 96

We propose a distributed-memory parallel GS (dmpGS) by implementing and 97 using a distributed-memory version of the sparse triangular Spike (stSpike) algorithm. 98 99 stSpike enables obtaining the solution of the system by solving independent sparse triangular subsystems and a smaller reduced triangular system. Solving this reduced 100 system constitutes a sequential computational bottleneck in dmpGS. The size of this 101 reduced system is equal to the number of nonzero columns in the lower off-diagonal 102blocks of the coefficient matrix. The computational cost of solving the reduced system 103 is proportional to its nonzero count. The communication volume of dmpGS is equal 104 105to the number of nonzero column segments in the off-diagonal blocks plus the reduced system size. Both of these communication and computational overheads highly depend 106on the sparsity structure of the coefficient matrix. 107

We note that solving the reduced system is embarrassingly parallel if the coeffi-108cient matrix is banded and diagonally dominant [50, 54]. In case the coefficient matrix 109 110 is not diagonally dominant, another way to alleviate the cost of solving the reduced system is to further parallelize the solution of the reduced system which has been done 111 iteratively [55], or recursively [15, 56]. Instead, we propose to minimize the size and 112 the nonzero count of the reduced system, together with the communication volume, 113and show that the resulting reduced system is so small that further parallelization of 114 the solution of the reduced system is often no longer needed. For attaining these min-115116 imization objectives, we propose a partitioning and reordering model that exploits the sparsity of the coefficient matrix. The proposed model consists of two phases. The 117 first phase is a row-wise partitioning of the coefficient matrix, whereas the second 118 phase is a row reordering within the row blocks induced by the partition obtained in 119 120the first phase.

For the first phase, we propose a novel hypergraph model that extends and enhances the conventional column-net model for simultaneously decreasing the reduced system size and the communication volume. We introduce vertex fixing, net anchoring and net splitting schemes within the recursive bipartitioning framework to encode the minimization of the number of nonzero column segments in the lower triangular part of the resulting partition.

For the second phase, we propose an intelligent in-block row reordering method with the aim of decreasing the computational costs of both forming the coefficient matrix of the reduced system once and solving the reduced system at each iteration. The rest of the paper is organized as follows. Section 2 provides the background

information on hypergraph and sparse matrix partitioning, and stSpike. In section 3,
we discuss the dmpGS algorithm along with its communication and computational
costs. The proposed partitioning and reordering model for dmpGS is introduced in
section 4. We provide the experimental results in section 5 and conclude in section 6.

135 **2. Background.**

2.1. Hypergraph Partitioning. A hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{N})$ consists of a set of 136 vertices $\mathcal{V} = \{v_i\}_{1 \leq i \leq n}$ and a set of nets $\mathcal{N} = \{n_j\}_{1 \leq j \leq m}$. Each net $n_j \in \mathcal{N}$ connects 137a subset of vertices in \mathcal{V} , which is referred to as the *pins* of n_i , and denoted by 138 $Pins(n_i, \mathcal{H})$. Each vertex v_i is assigned a weight $w(v_i)$ and each net n_i is assigned a 139cost $cost(n_i)$. $\Pi = \{\mathcal{V}_1, \mathcal{V}_2, \dots, \mathcal{V}_k\}$ denotes a K-way partition of \mathcal{H} , where parts are 140 mutually disjoint and exhaustive. The weight of a part is the sum of the weights of 141 vertices in that part. For a given partition, if a net connects at least one vertex in 142 a part, it is said to connect that part. Connectivity $\lambda(n_i)$ of net n_i is the number 143of parts connected by n_i . If a net n_i connects multiple parts (i.e. $\lambda(n_i) > 1$), it is 144called *cut*; and otherwise *internal* (i.e. $\lambda(n_i) = 1$). The set of cut nets is denoted by 145 \mathcal{N}_{cut} . The *cutsize* of Π is defined in various ways. Two most commonly used cutsize 146147 definitions are the *cut-net* and the *connectivity* metrics [18], which are respectively defined as 148

149 (2.1)
$$cs_{cutn}(\Pi) = \sum_{n \in \mathcal{N}_{cut}} cost(n), \text{ and } cs_{conn}(\Pi) = \sum_{n \in \mathcal{N}_{cut}} (\lambda(n) - 1)cost(n).$$

150Hypergraph partitioning (HP) is the problem of finding a K-way partition which minimizes the cutsize and satisfies the balance criterion $W_{max} \leq W_{avg}(1+\epsilon)$. Here, 151 ϵ is the given maximum allowable imbalance ratio; and W_{max} and W_{avg} respectively 152denote the maximum and average part weights. HP with fixed vertices ensures to 153assign some preassigned vertices which are called *fixed vertices* to the respective parts. 154The recursive bipartitioning (RB) is a widely used paradigm to obtain a K-way 155HP. It first partitions the hypergraph into two and then each part is further biparti-156 157tioned recursively until reaching the desired number of parts K. In order to encode the cut-net and connectivity metrics, cut-net removal and cut-net splitting methods 158are utilized in the RB-based HP, respectively [18]. 159

160 **2.2.** Sparse Matrix Partitioning with HP. Several HP models and methods 161 have been proposed and successfully utilized for obtaining matrix partitioning [8, 14, 162 19, 20, 25, 39, 65, 68]. Among these, the most relevant one is the *column-net* model 163 [18] that represents a given sparse matrix A as a hypergraph $\mathcal{H}_{CN}(A)$ in which nets 164 and vertices respectively represent columns and rows. In this model, vertex v_i is 165 added to the pin list of net n_j for each nonzero A(i, j) in A. Throughout the paper, 166 r_i and c_j respectively denote row i and column j.

A K-way ordered partition $\Pi = \langle \mathcal{V}_1, \mathcal{V}_2, \dots, \mathcal{V}_K \rangle$ of the column-net model $\mathcal{H}_{CN}(A)$ 167 is decoded as a partial reordering of the rows of A in such a way that the rows 168 corresponding to vertices in \mathcal{V}_k are ordered before the rows corresponding to the 169 vertices in \mathcal{V}_{ℓ} for $k < \ell$. This is a partial reordering since the rows corresponding to the 170 vertices in the same part can be ordered arbitrarily. Let \mathcal{B}_k^r denote the k^{th} row block 171which contains the rows corresponding to the vertices in \mathcal{V}_k . We consider a symmetric 172row-column reordering that yields a 2D grid structure of A. The submatrix consisting 173of the rows of \mathcal{B}_k^r and columns of ℓ^{th} column block \mathcal{B}_ℓ^c is referred as block- (k,ℓ) of 174A. A column is said to *connect* a row block \mathcal{B}_k^r if it contains at least one nonzero in 175176 \mathcal{B}_{k}^{r} . A column is called *cut* if it connects more than one row block. For a matrix with nonzero diagonal entries, each column connects a diagonal block and becomes a cut 177 178column if it connects at least one off-diagonal block.

In the column-net model with unit net cost, the partitioning objective using the connectivity and cut-net metrics (2.1) respectively encode the minimization of the number of nonzero column segments in off-diagonal blocks and the number of cut columns. The former partitioning objective is successfully utilized in encoding the 183 minimization of the row parallel SpMV operations [18].

184 2.3. Sparse Triangular Spike (stSpike) Algorithm. We describe stSpike for
 185 lower triangular systems since the algorithm for the upper triangular case is similar.
 186 Given a lower triangular linear system of equations

187 (2.2)
$$Ly = b,$$

a DS factorization of sparse lower triangular matrix L is computed as L=DS, where *D* is the lower block diagonal of L and S is the Spike matrix. These blocks are assumed to be obtained by matrix partitioning. Multiplying both sides of (2.2) from the left by D^{-1} , we obtain a modified system

192 (2.3)
$$Sy = g,$$

where $g = D^{-1}b$ and $S = D^{-1}L$. By splitting L = D + R, we obtain S = I + G where $G = D^{-1}R$, and R is the block off-diagonal part of L. The sparse triangular system DG = R with multiple right hand side vectors can be solved for the block rows of Gindependently with perfect parallelism.

The nonzero column segments of R constitute dense column segments (called spikes) in the off-diagonal blocks of S. The block diagonal of S is identity. Additional nonzeros (fill-in) are introduced within the off-diagonal blocks of S only in the locations below the top nonzero (having the smallest row index) for each nonzero column segment of R. The submatrix consisting of rows and columns C of S, namely $\hat{S} = S(C, C)$, constitutes an independent reduced system where C is the set of nonzero columns of R, i.e., cut columns of L. Then the reduced system is of the form

204 (2.4)
$$\widehat{S}\widehat{y} = \widehat{g},$$

where $\hat{g} = g(\mathcal{C})$ and $\hat{y} = y(\mathcal{C})$, which can be solved independent from the rest of the unknowns in y. After solving the reduced system, the only remaining computation for retrieving the solution of the original system is

208 (2.5)
$$y = g - D^{-1}(\hat{R}\hat{y}),$$

which can be obtained in perfect parallelism where $\widehat{R} = R(:, \mathcal{C})$ (in MATLAB notation). We only partially compute S just to form \widehat{S} , since forming S explicitly is expensive and requires a large amount of memory. Partial computation of S constitutes the factorization phase, whereas computation of \widehat{g} , solving (2.4) and (2.5) constitutes the solution phase of stSpike.

An example L matrix and the corresponding S and \widehat{S} matrices are shown in 214 215Figure 1. The reduced system indices $\mathcal{C} = \{1, 3, 4, 6, 7, 9, 11\}$ are colored in red and circled. The nonzeros that constitute the reduced system are bold and colored in red. 216217 The background colors of the original nonzeros and possible fill-in are green and blue, respectively. Depending on the sparsity structure of the corresponding column and 218block diagonal, spikes may not fill the entire column segment. For example, nonzero 219L(4,1) in block-(2,1) of L leads to the spike consisting of three nonzeros in the first 220 221 column of block-(2,1) of S.

3. Distributed-Memory Parallel GS (dmpGS) Algorithm. The pseudocode of dmpGS is given in Algorithm 3.1 for processor P_k in a K-processor system. Matrix A is assumed to be partitioned into K row blocks, where m_k denotes the



Fig. 1: Sparsity structure of L and resulting S and \widehat{S} matrices derived from stSpike.

number of rows in the k^{th} row block. In the algorithm, R_k , D_k and U_k respectively 225denote the k^{th} row block of the strictly block lower triangular, lower triangular part of 226 the block diagonal, and strictly upper triangular parts of A as shown in Figure 2. The number of columns in R_k , D_k and U_k are respectively $\sum_{i=1}^{k-1} m_i$, m_k and $\sum_{i=k}^{K} m_i$. 227228 f_k, g_k, x_k, h_k, w_k and z_k denote the local subvectors of size m_k that are computed by 229230 P_k . These subvectors are partitioned conformably with row-wise partitioning of A as shown in Figure 2. \widehat{S} , \widehat{x} and \widehat{g} respectively denote the $|\mathcal{C}| \times |\mathcal{C}|$ coefficient matrix, 231 $|\mathcal{C}| \times 1$ unknown and $|\mathcal{C}| \times 1$ right hand side vectors of the reduced system in stSpike. 232 C_k denotes the subset of C corresponding to the row indices in R_k . 233

In Algorithm 3.1, lines 2-7 denote the factorization phase of stSpike which com-234putes S. This phase is done only once after which we proceed with the GS iterations in 235236 lines 8-22. Each dmpGS iteration involves two SpMVs at lines 11 and 20, two vector subtraction operations at lines 12 and 22, an independent sparse triangular solve at 237 line 13, a reduced system solution at line 17, which enables independent sparse trian-238 gular solves at line 21. The upper and lower triangular SpMV operations are incurred 239 by the GS and stSpike algorithms, respectively. These two SpMV operations incur 240 241 communication of x-vector entries depending on the sparsity structures of the upper triangular U and lower triangular L matrices, respectively. Conformable partitioning 242 of the vectors avoids communication during vector subtraction operations. 243

At lines 9–10, communication operations are performed for local SpMV (line 24411). After P_k receives all necessary non-local x-vector entries, it forms its augmented 245vector \breve{x} . Each processor sends the selected entries of its g_k vector to P_1 (line 14) to 246247 form the right hand side vector \hat{q} (line 16) for the sequential solution of the reduced system to obtain \hat{x} (line 17). Here \hat{x} corresponds to those unknowns in x which are at 248the interface of the partitioning of L and obtaining them decouples the global lower 249 triangular system into independent much smaller systems. P_1 sends only those x-250vector entries that are required by other processors (line 18) so that each processor 251



Fig. 2: Four-way row-wise partition of matrix A and vectors x and f

Algorithm 3.1 Distributed-Memory Parallel Gauss Seidel (dmpGS) for processor P_k

In	nput: Submatrices R_k, D_k, U_k , and right-hand s	side subvector f_k
01	Dutput: Subvector x_k	
1: Ch	hoose an initial guess for x_k	
2: if	$2 \le k \le K - 1$ then	
3:	$G_k \leftarrow D_k^{-1} R_k$ \triangleright local partial sparse tria	ngular solve with multiple RHS
4:	Form and send \widehat{G}_k to processor P_1	
5: if	k = 1 then	
6:	Receive \widehat{G}_{ℓ} from P_{ℓ} for $2 \leq \ell \leq K-1$ to form \widehat{G}_{ℓ}	Ŷ
7:	$\widehat{S} \leftarrow \widehat{G} + I$	
8: wh	hile not converged do	
9:	Send required local x_k entries to respective pr	cocessors in $\{P_1, \ldots, P_{k-1}\}$
10:	Receive non-local x_{ℓ} entries from processors in	$\{P_{k+1},\ldots,P_K\}$ to form \breve{x}_k
11:	$h_k \leftarrow U_k \breve{x}_k$	$ ightarrow m local \ SpMV$
12:	$h_k \leftarrow f_k - h_k$	
13:	$g_k \leftarrow D_k^{-1} h_k$	\triangleright local sparse triangular solve
14:	if $2 \le k \le K-1$ then Send $\{g_k(i)\}_{i \in \mathcal{C}_k}$ to proc	tessor P_1
15:	if $k = 1$ then	
16:	Receive $\{g_{\ell}(i)\}_{i \in \mathcal{C}_k}$ from P_{ℓ} for $2 \leq \ell \leq K$ –	- 1 to form \widehat{g}
17:	$\widehat{x} \leftarrow \widehat{S}^{-1}\widehat{g}$	\triangleright solve reduced system
18:	Send \hat{x} entries to requiring processors	
19:	if $k \neq 1$ then Receive required \hat{x} -entries to for	$\operatorname{rm} \bar{x}_k$
20:	$z_k \leftarrow R_k \bar{x}_k$	$ ightarrow m local \ SpMV$
21:	$w_k \leftarrow D_k^{-1} z_k$	\triangleright local sparse triangular solve
22:	$x_k \leftarrow g_k - w_k$	

252 P_k forms its \bar{x} vector (line 19) to perform local SpMV (line 20).

The communication overhead in each iteration of dmpGS is as follows. The communication volume incurred by h = Ux (line 11) and z = Rx (line 20) are equal to the number of nonzero column segments in the off-diagonal blocks of U and L, respectively. Thus the communication volume required by these two SpMVs is equal to the total number of off-diagonal nonzero column segments in A (offD_nzCol_seg(A)). The volume of communication incurred at line 16 is equal to the size of the reduced system, |C|. Therefore, the total communication volume of dmpGS is

260 (3.1) $\operatorname{commVol} = \operatorname{offD_nzCol_seg}(A) + |\mathcal{C}|.$

Note that the different row blocks (R_k) seem to vary in the number of columns because of the triangular structure of the problem. On the other hand, the computational load imbalance is alleviated by the proposed partitioning model which also gathers most of the nonzeros to the diagonal blocks.

4. The Proposed Partitioning and Reordering Model. We propose a twophase model for reducing the communication overhead of dmpGS while maintaining computational balance as well as reducing the sequential computational overhead incurred by solving the reduced system at each iteration. This computational overhead is proportional to the number of nonzeros in the off-diagonals of \hat{S} . In subsection 4.1, we propose a novel HP model as the first phase which simultaneously encodes the minimization of the reduced system size $|\mathcal{C}|$ and the communication volume. Decreasing

 $|\mathcal{C}|$ is important not only because it directly contributes to reducing the commu-272 273nication volume, but it also relates to decreasing the computational overhead. In subsection 4.2, we propose an in-block reordering method as the second phase which 274refines the improvement further by decreasing the number of nonzeros in S. We pro-275vide the illustrations showing the effect of the proposed partitioning and reordering 276model on a sample matrix in subsection 4.3. 277

4.1. Hypergraph Partitioning Model. The partitioning objective in this 278phase is minimizing the sum of communication volume overhead (3.1) and sequential 279overhead costs with proper scaling: 280

281

8

281
$$PartObj = commVol + (\alpha - 1)|\mathcal{C}|$$

282
$$= (offD_nzCol_segs(A) + |\mathcal{C}|) + (\alpha - 1)|\mathcal{C}|$$

$$(\texttt{offD_nzCol_segs}(A) + |c|) + (\alpha)$$

$$= \texttt{offD_nzCol_segs}(A) + \alpha |\mathcal{C}|$$

Here α denotes the scaling factor between the effect of the reduced system size and 285the number of off-diagonal nonzero column segments on the overall overhead. 286

4.1.1. Definitions and Layout. We define a column as *L*-cut if it connects at 287 least one off-diagonal block in the lower triangular part. That is, a column c_i in k^{th} 288 column block \mathcal{B}_k^c is L-cut if it connects a row block \mathcal{B}_ℓ^r with $\ell > k$. Since L-cut columns 289of A are the nonzero columns of R, the number of L-cut columns $(L-cut_cols(A))$ is 290 291 equal to the reduced system size, $|\mathcal{C}|$. Therefore, the partitioning objective (4.1) can 292 be rewritten as

293 (4.2)
$$PartObj = offD_nzCol_segs(A) + \alpha(L-cut_cols(A)).$$

Let $\mathcal{H}_{CN}(A) = (\mathcal{V}, \mathcal{N})$ be the column-net hypergraph of an $m \times m$ sparse matrix A 294 with nonzero diagonal entries. An ordered partition $\Pi_K = \langle \mathcal{V}_1, \mathcal{V}_2, \dots, \mathcal{V}_K \rangle$ of $\mathcal{H}_{CN}(A)$ 295is decoded as a partial symmetric row and column reordering of A as explained in 296section 2.2. Each net n_i of $\mathcal{H}_{CN}(A)$ connects vertex v_i since $A(i,i) \neq 0$ for each 297 $1 \leq i \leq m$. A net n_i with $v_i \in \mathcal{V}_k$ is called *L*-cut if it connects at least one vertex part 298 \mathcal{V}_{ℓ} such that $\ell > k$. The set of L-cut nets is denoted as \mathcal{N}_{Lcut} . We define a new type 299 300 of cutsize, which we call the *L*-cut-net metric, as

301 (4.3)
$$cs_{Lcut}(\Pi_K) = \sum_{n \in \mathcal{N}_{Lcut}} cost(n).$$

Finally, the cost of partition Π_K is defined as the sum of connectivity metric with 302 unit net cost and L-cut-net metric with net cost α , i.e., 303

304 (4.4)
$$cost_{conn+Lcut}(\Pi_K) = \sum_{n \in \mathcal{N}_{cut}} (\lambda(n)-1) + \alpha |\mathcal{N}_{Lcut}|.$$

Here, each cut net n incurs $\lambda(n) - 1$, and each L-cut net incurs α to the cutsize. 305

LEMMA 4.1. A column c_i of A is L-cut iff net n_i of $\mathcal{H}_{CN}(A)$ is L-cut. 306

Proof. Due to symmetric row-column ordering, c_i is in \mathcal{B}_k^c iff r_i is in \mathcal{B}_k^r , which 307 corresponds to $v_i \in \mathcal{V}_k$. Furthermore, c_i connects \mathcal{B}_{ℓ}^r iff n_i connects \mathcal{V}_{ℓ} . Therefore, c_i 308 in \mathcal{B}_k^c connects \mathcal{B}_ℓ^r iff n_i with $v_i \in \mathcal{V}_k$ connects \mathcal{V}_ℓ , where $\ell > k$. Π 309

PROPOSITION 4.2. Minimizing $cost_{conn+Lcut}(\Pi_K)$ for a K-way partition Π_K of 310 311 $\mathcal{H}_{CN}(A)$ corresponds to minimizing the partitioning objective (4.2).

9

312 Proof. By Lemma 4.1, the number of L-cut nets in $\mathcal{H}_{CN}(A)$ is equal to the 313 number of L-cut columns in A. Thus $\alpha |\mathcal{N}_{Lcut}| = \alpha(L-\text{cut_cols}(A))$. Furthermore, 314 it is known by [18] that $\sum_{n \in \mathcal{N}_{cut}} (\lambda(n)-1) = \text{offD_nzCol_segs}(A)$.

Each vertex is associated with a weight equal to the number of nonzeros in the respective row of the matrix, i.e., $w(v_i) = nnz(A(i,:))$. Thus, the partitioning constraint of maintaining balance on part weights approximately encodes the computational load balance during aggregate two triangular SpMVs (lines 11 and 20) and two triangular solves (lines 13 and 21).

The cut-net splitting technique has been successfully used within the RB framework to encode the minimization of the connectivity metric [18]. However to the best of our knowledge, there exists no tool or model for encoding the minimization of the *L*-cut-net metric in the literature. We propose to use the RB framework with novel net anchoring and splitting schemes to encode the minimization of the *L*-cut-net metric.

4.1.2. Recursive Bipartitioning Model. At each RB step, an ordered bipartition $\Pi_2 = \langle \mathcal{V}_U, \mathcal{V}_L \rangle$ of \mathcal{V} is decoded as ordering the vertices of \mathcal{V}_L after those of \mathcal{V}_U . Here \mathcal{V}_U and \mathcal{V}_L denote the upper and lower vertex parts, respectively. In RB, the concept of *L*-cut net takes a special form. In a bipartition $\Pi_2 = \langle \mathcal{V}_U, \mathcal{V}_L \rangle$, a net n_i is *L*-cut if v_i is assigned to \mathcal{V}_U and n_i connects at least one vertex v_j such that $v_j \in \mathcal{V}_L$. The partitioning objective at each RB step is to minimize

332 (4.5)
$$cost_{RB}(\Pi_2) = |\mathcal{N}_{cut}| + \alpha |\mathcal{N}_{Lcut}|.$$

For encapsulating the connectivity and L-cut net metrics simultaneously, each 333 net n_i in $\mathcal{H}_{CN}(A)$ is replicated as two different kinds of nets, namely conn-net n_i^c 334 and lcn-net n_i^{ℓ} . Here, conn-nets encapsulate the connectivity metric whereas lcn-335 nets encapsulate the L-cut-net metric. The motivation for net replication is the 336 337 requirement of different net splitting and net removal procedures for encoding the connectivity and L-cut-net metrics at each RB step. In order to encapsulate the RB 338 339 objective (4.5), we assign unit cost to the conn-nets and cost α to the lcn-nets. We refer to the hypergraph formed by these replicated nets as \mathcal{H} . 340

We extend $\mathcal{H} = (\mathcal{V}, \mathcal{N})$ into a hypergraph $\mathcal{H}' = (\mathcal{V}', \mathcal{N}')$ so that minimizing the number of conventional cut nets in \mathcal{H}' encodes minimizing (4.5). We introduce new fixed vertices $v_U \in \mathcal{V}_U$ and $v_L \in \mathcal{V}_L$ to form the extended vertex set $\mathcal{V}' = \mathcal{V} \cup \{v_U, v_L\}$. We represent each lcn-net n_i^{ℓ} in \mathcal{H} as a pair of nets \hat{n}_i^{ℓ} and \check{n}_i^{ℓ} in \mathcal{H}' . \hat{n}_i^{ℓ} is same as n_i^{ℓ} except it is *U*-anchored (connects v_U). \check{n}_i^{ℓ} is a 2-pin *L*-anchored net which connects v_L and v_i . That is, for each net n_i in $\mathcal{H}_{CN}(A)$, \mathcal{H}' contains nets n_i^c , \hat{n}_i^{ℓ} and \check{n}_i^{ℓ} , where

347 $Pins(n_i^c, \mathcal{H}') = Pins(n_i, \mathcal{H}_{CN}(A)),$

348 $Pins(\hat{n}_i^{\ell}, \mathcal{H}') = Pins(n_i, \mathcal{H}_{CN}(A)) \cup \{v_U\} \text{ and }$

 $\frac{349}{350} \quad Pins(\check{n}_i^\ell, \mathcal{H}') = \{v_i, v_L\}.$

The nets in the extended hypergraph for a sample 3-pin net are shown in Figure 3.

We form \mathcal{H}' at the beginning and apply RB steps until reaching the desired part count, K. The resulting K-way partition Π'_K of \mathcal{H}' induces a K-way partition Π_K of $\mathcal{H}_{CN}(A)$. \mathcal{H} is an in-between hypergraph introduced for the sake of clarity of presentation and is not constructed during implementation. We explain the proposed net splitting and removal methods on \mathcal{H} , and show the correspondence on \mathcal{H}' . We consider that each bipartition $\Pi'_2 = \langle \mathcal{V}'_U, \mathcal{V}'_L \rangle$ of \mathcal{H}' induces a bipartition $\Pi_2 = \langle \mathcal{V}_U, \mathcal{V}_L \rangle$



Fig. 3: Net n_i in $\mathcal{H}_{CN}(A)$ is replicated as conn-net n_i^c and lcn-net n_i^ℓ to form \mathcal{H} . Net n_i^ℓ in \mathcal{H} is represented by a pair of nets \hat{n}_i^ℓ and \check{n}_i^ℓ in \mathcal{H}' .

of \mathcal{H} . Here \mathcal{H} and \mathcal{H}' refer to the respective hypergraphs just before the current RB step. New hypergraphs \mathcal{H}_U and \mathcal{H}_L are constructed according to $\Pi_2 = \langle \mathcal{V}_U, \mathcal{V}_L \rangle$ as follows. For both conn- and lcn-nets, each internal net in \mathcal{V}_L and \mathcal{V}_U is respectively included in \mathcal{N}_L and \mathcal{N}_U as is. In the net splittings, a new conn- or lcn-net is added to the net list of \mathcal{H}_U or \mathcal{H}_L only if it has more than one pin. The single-pin nets are discarded since they cannot contribute to the cutsize in the following RB steps.

For cut conn-nets, we apply the conventional cut-net splitting procedure [18] to encapsulate the connectivity metric. If a conn-net n_i^c is cut, then n_i^c is split into two pin-wise disjoint nets in \mathcal{H}_U and \mathcal{H}_L such that

367 $Pins(n_i^c, \mathcal{H}_U) = Pins(n_i^c, \mathcal{H}) \cap \mathcal{V}_U$, and $Pins(n_i^c, \mathcal{H}_L) = Pins(n_i^c, \mathcal{H}) \cap \mathcal{V}_L$.

For lcn-nets, we introduce a hybrid cut-net splitting/removal method in order to correctly encapsulate the *L*-cut-net metric. At each RB step, for each net pair $(\hat{n}_i^{\ell}, \check{n}_i^{\ell})$ in a bipartition Π' , we consider the state of n_i^{ℓ} in Π where $Pins(n_i^{\ell}, \mathcal{H}) =$ $Pins(\hat{n}_i^{\ell}, \mathcal{H}') - \{v_U\}$ for ease of understanding. If an lcn-net n_i^{ℓ} is not internal, then it can be *L*-cut or "cut but not *L*-cut".

If n_i^{ℓ} is L-cut, then we apply cut-net removal for n_i . This is because when n_i is 373 L-cut in an RB step, it also becomes L-cut in the final K-way partition. Hence there 374 is no need to track this net anymore and we do not include it in further bipartitions. 375 If n_i^{ℓ} is cut but not L-cut, then we apply net removal towards \mathcal{H}_U and net-L-376 splitting towards \mathcal{H}_L . That is, n_i^{ℓ} is added to \mathcal{H}_L as $Pins(n_i^{\ell}, \mathcal{H}_L) = Pins(n_i^{\ell}, \mathcal{H}) \cap \mathcal{V}_L$. 377 This is because n_i^{ℓ} cannot be L-cut in further bipartitionings of \mathcal{H}_U but it has the 378 potential of becoming L-cut in further bipartitionings of \mathcal{H}_L . In the extended hyper-379 graph context, this corresponds to adding lcn-net pair $(\hat{n}_i^{\ell}, \check{n}_i^{\ell})$ to \mathcal{H}'_L such that 380

381
$$Pins(\hat{n}_i^{\ell}, \mathcal{H}_L') = (Pins(n_i^{\ell}, \mathcal{H}') \cap \mathcal{V}_L') \cup \{v_U\}, \text{ and } Pins(\check{n}_i^{\ell}, \mathcal{H}_L') = \{v_i, v_L\}.$$

Figure 4 shows all possible cases for a sample lcn-net. The first, second, third and last horizontal layers respectively show the bipartition Π'_2 of \mathcal{H}' ; the corresponding bipartition Π_2 of \mathcal{H} ; \mathcal{H}_U and \mathcal{H}_L induced by Π_2 ; and the corresponding \mathcal{H}'_U and \mathcal{H}'_L induced by Π'_2 . If n_i^{ℓ} is *L*-cut in \mathcal{H} as in Figure 4a, both \hat{n}_i^{ℓ} and \check{n}_i^{ℓ} are cut in Π'_2 . If n_i^{ℓ} is cut but not *L*-cut as in Figure 4b, or if n_i^{ℓ} is internal to \mathcal{V}_L as in Figure 4c, then only \hat{n}_i^{ℓ} is cut. Otherwise, if n_i^{ℓ} is internal to \mathcal{V}_U as in Figure 4d, then only \check{n}_i^{ℓ} is cut.

LEMMA 4.3. Consider the bipartition $\Pi_2 = \langle \mathcal{V}_U, \mathcal{V}_L \rangle$ of \mathcal{H} induced by a bipartition $\Pi'_2 = \langle \mathcal{V}'_U, \mathcal{V}'_L \rangle$ of \mathcal{H}' in an RB step. If a net is L-cut in Π_2 , then it incurs 2 cut nets in Π'_2 . Conversely, if a net is not L-cut in Π_2 , then it incurs 1 cut net in Π'_2 .

391 Proof. If n_i^{ℓ} is L-cut in Π_2 of \mathcal{H} , then $v_i \in \mathcal{V}_U$ and n_i^{ℓ} connects a vertex v_j such 392 that $v_j \in \mathcal{V}_L$. In Π'_2 of \mathcal{H}' , \hat{n}_i^{ℓ} is cut since it connects $v_i \in \mathcal{V}'_U$ and $v_j \in \mathcal{V}'_L$; and \check{n}_i^{ℓ} is 393 also cut since it connects $v_i \in \mathcal{V}'_U$ and $v_L \in \mathcal{V}'_L$.



(a) n_i^{ℓ} is *L*-cut (b) n_i^{ℓ} is cut, not *L*-cut (c) n_i^{ℓ} is internal to \mathcal{V}_L (d) n_i^{ℓ} is internal to \mathcal{V}_U Fig. 4: All cases for n_i^{ℓ} and corresponding net pair $(\hat{n}_i^{\ell}, \check{n}_i^{\ell})$ after bipartition $\langle \mathcal{V}'_U, \mathcal{V}'_L \rangle$.

If n_i^{ℓ} is not *L*-cut and $v_i \in \mathcal{V}_L$ in Π_2 , then \hat{n}_i^{ℓ} is cut in Π'_2 because it connects $v_U \in \mathcal{V}'_U$ and $v_i \in \mathcal{V}'_L$; but \check{n}_i^{ℓ} is not cut since both v_i and v_L are in \mathcal{V}'_L .

396 If n_i^{ℓ} is not *L*-cut and $v_i \in \mathcal{V}_U$ in Π_2 , then n_i^{ℓ} should be internal to \mathcal{V}_U , because 397 otherwise any pin in \mathcal{V}_L would make n_i^{ℓ} to be *L*-cut. In Π'_2 , net \hat{n}_i^{ℓ} is internal to \mathcal{V}'_U 398 since both v_i and v_U are in \mathcal{V}'_U ; but \tilde{n}_i^{ℓ} is cut since it connects $v_i \in \mathcal{V}'_U$ and $v_L \in \mathcal{V}'_L$. \square

PROPOSITION 4.4. Minimizing the conventional cut-net metric for the bipartition Π'_2 of \mathcal{H}' encodes minimizing $cost_{RB}(\Pi_2)$ defined in (4.5).

401 Proof. By Lemma 4.3, each L-cut net in Π_2 incurs 2 cut nets in Π'_2 , whereas 402 all remaining nets in Π_2 incur 1 cut net in Π'_2 . Since the cost of lcn-nets is α , 403 the cutsize incurred by lcn-nets in Π'_2 is $\alpha(|\mathcal{N}_{Lcut}| + |\mathcal{N}|)$. Since conn-nets are of 404 unit cost, they incur $|\mathcal{N}_{cut}|$ to the cutsize of Π'_2 . Hence the total cutsize of Π'_2 is 405 $|\mathcal{N}_{cut}| + \alpha |\mathcal{N}_{Lcut}| + \alpha |\mathcal{N}|$. Since $\alpha |\mathcal{N}|$ is constant, minimizing the cutsize of Π'_2 is 406 equivalent to minimizing $|\mathcal{N}_{cut}| + \alpha |\mathcal{N}_{Lcut}|$, which is $cost_{RB}(\Pi_2)$.

Figure 5 shows an example 2-level RB in terms of lcn-nets in \mathcal{H} and the corresponding 4-way matrix partitioning. The *L*-cut nets n_1^{ℓ} , n_2^{ℓ} and n_6^{ℓ} and the corresponding *L*-cut columns c_1 , c_2 , and c_6 of *A* are colored in red background. n_2^{ℓ} is *L*-cut in the first level RB and discarded in the future bipartitions. This is because column c_2 is already counted as *L*-cut due to nonzero A(6, 2) and should not be counted as *L*-cut again due to nonzero A(4, 2) in further bipartitions.

⁴¹³ Note that the *L*-cut net definition can be considered to be similar to the left-cut ⁴¹⁴ net defined in [1] for encapsulating the profile minimization, but the net splitting and



Fig. 5: Sample 2-level RB showing lcn-nets and corresponding matrix partitioning.

415 removal strategies are quite different for encapsulating the objective of our problem.

⁴¹⁶ THEOREM 4.5. Recursively bipartitioning \mathcal{H}' by minimizing the cutsize according ⁴¹⁷ to the cut-net metric and applying the proposed net splitting and removal strategies ⁴¹⁸ until reaching K parts encode minimizing the partitioning objective (4.2).

Proof. By Proposition 4.4, recursively bipartitioning \mathcal{H}' by minimizing the con-419 ventional cut-net metric encodes minimizing $cost_{RB}(\Pi_2)$ at each RB step. We show 420 that this encodes minimizing $cost_{conn+Lcut}(\Pi_K)$. Proposed net splitting and removal 421 strategies ensure that an L-cut net in Π_K is also L-cut in Π_2 in exactly one RB step. 422 Since an L-cut net contributes α to both $cost_{RB}(\Pi_2)$ and $cost_{conn+Lcut}(\Pi_K)$, minimiz-423 ing $\alpha |\mathcal{N}_{Lcut}|$ in each bipartition Π_2 encodes minimizing $\alpha |\mathcal{N}_{Lcut}|$ in Π_K . Furthermore, 424minimizing the number of cut nets $|\mathcal{N}_{cut}|$ at each RB step and applying the cut-net 425splitting procedure encodes minimizing the connectivity metric $\sum_{n \in \mathcal{N}_{cut}} (\lambda(n)-1)$ [18]. Therefore, minimizing the cutsize for each bipartition Π'_2 of \mathcal{H}' encodes minimizing 426 427 $cost_{conn+Lcut}(\Pi_K)$; hence by Proposition 4.2, this corresponds to the partitioning 428 objective (4.2). 429

430 **4.2. Reordering within Row Blocks.** Consider the *K*-way block structure 431 (e.g. Figure 2) of *A* induced by the partial symmetric row-column permutation ob-432 tained by the HP model (section 4.1). We perform row reordering within the k^{th} 433 row block of *A* by considering nonzeros of the k^{th} row block R_k of *R*. The resulting 434 row reordering within the k^{th} row block of *A* is symmetrically applied to the columns 435 of the k^{th} column block of *A*. R_k is an $m_k \times z_k$ matrix where $z_k = \sum_{i=1}^{k-1} m_i$. For 436 simplicity, we assume a local indexing for the rows of R_k so that R_k consists of rows 437 r_i with $1 \le i \le m_k$.

438 Recall that in stSpike, fill-in may arise below the top nonzero of each spike in R_k . 439 The top nonzero of a spike c_j in R_k is the nonzero with the minimum row index, i.e., 440 $top(c_j, R_k) = \min\{i : R_k(i, j) \neq 0, 1 \leq i \leq m_k\}$. We define the *height* of a spike c_j in R_k 441 as the number of reduced system row indices between $top(c_j, R_k)$ and m_k inclusively, 442 i.e.,

443 (4.6)
$$height(c_i, R_k) = |\{i : top(c_i, R_k) \le i \le m_k, i \in \mathcal{C}_k\}|,$$

since only the rows with indices in C_k may contribute to the nonzero count of \widehat{S} . The height of a spike in R_k constitutes an upper bound on the nonzero count (including the fill-in) of the corresponding column in \widehat{S} . In Figure 1b, the heights of the spikes are as follows: $height(c_1, R_2)=3$, $height(c_3, R_2)=2$; $height(c_1, R_3)=1$, $height(c_4, R_3)=2$,

448 $height(c_6, R_3)=1$, and $height(c_7, R_3)=2$. The height of a non-spike column is assumed

449 to be zero. The objective of in-block reordering is to minimize the *total height*

450 (4.7)
$$\sum_{k=2}^{K-1} \sum_{j=1}^{z_k} height(c_j, R_k),$$

which constitutes an upper bound on the nonzero count in off-diagonal blocks of \widehat{S} . The last block R_K does not contribute nonzeros to \widehat{S} since \mathcal{C}_K is empty. Reordering within different blocks are completely independent and can be done concurrently.

One straightforward approach is placing the rows whose indices are not among 454 \mathcal{C}_k to the bottom of R_k to avoid the nonzeros of the rows that are not in \mathcal{C}_k to 455contribute to (4.7). Let $\overline{R}_k = R_k(\mathcal{C}_k, :)$ be the $|\mathcal{C}_k| \times z_k$ submatrix of R_k consisting of 456the rows with indices in C_k . Then the problem is reduced to reordering only those 457rows of \overline{R}_k since the rest of the rows at the bottom of R_k do not have an impact on 458(4.7). The reordering objective for each \overline{R}_k is to minimize $\sum_{j=1}^{z_k} height(c_j, \overline{R}_k)$, with 459a simplified height definition, $height(c_j, \overline{R}_k) = |\mathcal{C}_k| + 1 - top(c_j, \overline{R}_k)$. Then the total 460height minimization problem is formulated in general as: Given any sparse matrix 461 $H \in \mathbb{R}^{\ell \times n}$, find a row reordering P that minimizes $\sum_{j=1}^{n} (\ell + 1 - top(c_j, PH)).$ 462

463 THEOREM 4.6. The total height minimization problem (THMP) is NP-hard.

Proof. We reduce the profile minimization problem (PMP) [1, 46], which is known 464 to be NP-hard [26, 44], to THMP as follows. Given a symmetric matrix $V \in \mathbb{R}^{n \times n}$ with 465nonzero diagonal entries, the objective of PMP is finding a symmetric row/column reordering P that minimizes $\sum_{j=1}^{n} (j - top(c_j, PVP^T))$. This minimization objective is equivalent to maximizing $\sum_{j=1}^{n} top(c_j, PVP^T)$, since $\sum_{j=1}^{n} j$ is constant. Any instance PVP^T of PMP can be mapped to an instance PV of THMP by simply 466 467 468469 removing the column reordering as $(PVP^T)P = PV$. Note that the minimization 470 objective of THMP, which is $\sum_{j=1}^{n} (n+1-top(c_j, PV))$, is equivalent to maximizing 471 $\sum_{j=1}^{n} top(c_j, PV)$, since $\sum_{j=1}^{n} (n+1)$ is constant. Thus, PVP^T is a solution of PMP iff PV is a solution of THMP since the column reordering itself has no effect on 472 473 $\sum_{j=1}^{n} top(c_j, PVP^T) = \sum_{j=1}^{n} top(c_j, PV)$. If there had been a polynomial-time solu-474tion to THMP, then one could solve PMP in polynomial time by just applying the 475row reordering obtained by THMP on the columns as well. Therefore, PMP can be 476reduced to THMP in polynomial time; and since PMP is NP-hard, then so is THMP. 477

Algorithm 4.1 presents the pseudocode of the proposed heuristic for reordering 478 the rows of R_k . The efficient implementation of this algorithm requires accessing 479the nonzeros of both rows and columns of \overline{R}_k , so it is stored both in CSR and CSC 480 formats. $Cols(r_i)$ denotes the set of columns in row r_i , whereas $Rows(c_i)$ denotes 481 the set of rows in column c_i . Degree of a row or column is defined as the number of 482 nonzeros in that row or column, i.e., $deg(r_i) = |Cols(r_i)|$ and $deg(c_i) = |Rows(c_i)|$. In 483 lines 3-5, we define the *load* of each row r_i as the sum of degrees of columns c_i such 484 that $\overline{R}_k(i,j) \neq 0$. 485

The greedy choice utilized in the proposed heuristic is to order the rows with 486smaller degrees to upper positions of \overline{R}_k since placing denser rows to upper positions 487 incurs more height in (4.7). We further improve our greedy approach by using dynamic 488 489 row degrees during the row selection process. When a row is selected, the degree of each unselected row is decremented by the number of its nonzeros having the 490same column index with the nonzeros in the selected row. Since the nonzeros in a 491selected row already determine the heights of the respective columns, we do not need 492to consider the rest of the nonzeros of these columns in future row selections. When 493

Algorithm 4.1 Proposed in-block reordering for R_k where $2 \le k \le K-1$ **Input:** $R_k \in \mathbb{R}^{m_k \times z_k}$ and set of reduced-system row indices C_k of R_k . **Output:** the permutation vector *perm* of R_k . 1: Place the rows r_i with $i \notin C_k$ to the last $m_k - |C_k|$ indices in any order 2: Consider submatrix $\overline{R}_k = R_k(\mathcal{C}_k, :)$ of R_k consisting of rows r_i with $i \in \mathcal{C}_k$ for each row r_i of \overline{R}_k do 3: 4: $load(r_i) \leftarrow 0$ for each column $c_i \in Cols(r_i)$ do $load(r_i) \leftarrow load(r_i) + deg(c_i)$ 5: 6: for $d \leftarrow 0$ to max_row_deg do $\mathcal{S}(d) \leftarrow \{r_i : deg(r_i) = d\}$ 7: $indx \leftarrow 0$ while $indx < |\mathcal{C}_k|$ do 8: $d^* \leftarrow \min\{d : \mathcal{S}(d) \neq \emptyset\}$ 9: $r_{i^*} \leftarrow \operatorname{argmax}_{r_i \in \mathcal{S}(d^*)} load(r_i)$ \triangleright Select $r_{i^*} \in \mathcal{S}(d^*)$ with maximum load 10: 11: $indx \leftarrow indx + 1$ $perm(indx) \leftarrow r_{i^*}$ 12: $\mathcal{S}(d^*) \leftarrow \mathcal{S}(d^*) - \{r_{i^*}\}$ 13:for each column $c_i \in Cols(r_{i^*})$ do 14: $Rows(c_i) \leftarrow Rows(c_i) - \{r_{i^*}\}$ 15: $Cols(r_{i^*}) \leftarrow Cols(r_{i^*}) - \{c_i\}$ 16:for each row $r_{i'} \in Rows(c_i)$ do 17: $Cols(r_{i'}) \leftarrow Cols(r_{i'}) - \{c_i\}$ 18: $load(r_{i'}) \leftarrow load(r_{i'}) - deg(c_j)$ 19: $\mathcal{S}(deg(r_{i'})) \leftarrow \mathcal{S}(deg(r_{i'})) - \{r_{i'}\}$ 20: $deg(r_{i'}) \leftarrow deg(r_{i'}) - 1$ 21: 22: $\mathcal{S}(deg(r_{i'})) \leftarrow \mathcal{S}(deg(r_{i'})) \cup \{r_{i'}\}$

selecting a row among rows with the same degree, load values of the rows are used as a tie-breaking strategy. A row with a higher load is preferred to be selected since it will lead to a larger amount of decrease on the degrees of unselected rows.

In Algorithm 4.1, S(d) denotes the set of rows with degree d. Due to dynamic row degrees, at each iteration we find the minimum degree d^* (line 9). Then we choose the row r_{i^*} in $S(d^*)$ with the maximum load (line 10). After r_{i^*} is selected, all remaining nonzeros in each column c_j with $R_k(i^*, j) \neq 0$ are deleted as in lines 15-18. For each unselected row $r_{i'}$ with $R_k(i', j) \neq 0$, we dynamically update the load and degree of $r_{i'}$, and the respective degree sets (lines 19-22).

Recall that forming \hat{S} in dmpGS requires the computation of nonzeros up to the largest reduced system row index and any entry beyond that is not required to be computed for each row block. Hence the total height (4.7) also gives the computational cost of forming \hat{S} since we place \overline{R}_k at the top of R_k for each 1 < k < K.

4.3. Illustration. Figure 6 illustrates the effect of applying the proposed partitioning and reordering model for K=8 on a sample matrix (msc23052) from the SuiteSparse Matrix Collection [24]. The nonzero structure of the original matrix, the structure obtained after applying the proposed HP model and the final structure after the proposed in-block reordering are shown in order. Below each ordering of A, the resulting Spike matrix (S) is shown, including the nonzeros of the reduced system (\hat{S}) which are highlighted with red circles. As seen in the figure, the proposed partitioning and reordering model significantly reduces the nonzero count of the reduced system.



Fig. 6: Nonzero structure of msc23052: (a) before ordering, (b) after HP for K = 8, (c) after HP and in-block reordering; (d),(e),(f) the respective Spike (S) matrices (the reduced system (\hat{S}) nonzeros are circled in red color).

- 515 For example, the number of nonzeros in $\widehat{S}-I$ in Figures 6d, 6e, and 6f are 277,113,
- 516 3,593, and 811, respectively. Note that these numbers may seem to be much larger 517 than the ones appearing in the figures because of the overlapping red circles.

Notice that the off-diagonal blocks because of the overlapping red childes. Notice that the proposed HP model gathers most of the nonzeros to the diagonal blocks so that the off-diagonal blocks become very sparse. Then, the proposed in-block reordering method gathers the reduced-system nonzeros to the upper left corner of the respective off-diagonal block (Figure 6f). This is because we agglomerate the reduced system row indices to the top within each block, and we apply the resulting row reordering to the columns symmetrically. Within each off-diagonal block, gathering

the rows with reduced-system indices to the top corresponds to agglomerating the columns with these indices, which are actually all the columns having nonzeros, to the left. An exception is the first column block since no row reordering is performed for the first row block.

5. Experiments. We use the HSL software package MC64 [29] for scaling and permuting the coefficient matrices to avoid a singular *L*. We select the MC64 option that maximizes the product of the diagonal entries and then scales to make the absolute value of diagonal entries one and the off-diagonal entries less than or equal to one. For symmetric matrices, in order not to destroy the symmetry, we apply the symmetric MC64 if the main diagonal is already zero-free. Otherwise, we apply the nonsymmetric MC64 to obtain a zero-free main diagonal. For unsymmetric matrices, we just apply the nonsymmetric MC64.

Kind ID	Kind Name	Sym	Unsym	Total
1	structural	48	4	52
2	circuit simulation	2	46	48
3	economic	1	33	34
4	semiconductor device	0	33	33
5	computational fluid dynamics	6	27	33
6	2D/3D	19	9	28
7	power network	14	13	27
8	optimization	20	3	23
9	model reduction	13	3	16
10	chemical process simulation	0	15	15
11	theoretical/quantum chemistry	14	0	14
12	electromagnetics	6	4	10
13	thermal	5	4	9
14	materials	2	4	6
15	weighted graph	1	5	6
16	acoustics, oceanography, counter-ex., analytics	4	1	5
	All	155	204	359

Table 1: Number of instances among different matrix kinds in the dataset

536 The experiments are conducted on an extensive dataset obtained from the SuiteSparse Matrix Collection [24]. For sufficiently coarse-grained parallel processing, we select real square matrices that have more than 20,000 rows and between 100,000 538 and 20,000,000 nonzeros. There are 199 symmetric and 208 unsymmetric matrices in 539SuiteSparse satisfying these properties at the time of experimentation. 44 symmetric 540and 4 unsymmetric matrices are eliminated because they are singular. The remaining 541are 155 symmetric and 204 unsymmetric, a total of 359 sparse matrices on which we 542 conduct experiments. Table 1 shows the number of instances for each matrix kind. 543Kinds are sorted in decreasing order of instance count. The kinds having less than 5445 instances in our dataset (acoustics, chemical oceanography, counter-example and 545data analytics) are grouped as one kind. 546

5.1. Partitioning Quality. We tested the performance of the proposed parti-547548 tioning algorithm described in subsection 4.1 against the partitioning quality of the conventional column-net HP with connectivity metric (cnHP) and graph partitioning 549(GP) models. For both cnHP and GP, vertex weights are set as the number of nonze-550 ros in the respective rows whereas nets and edges are assigned unit cost. In cnHP, the objective is to minimize the number of nonzero off-diagonal column segments. In 552GP, the objective is to minimize the number of nonzeros in the off-diagonal blocks. 553For unsymmetric matrices, GP is applied on $|A| + |A^T|$. The well-known partitioning 554tools METIS [40] and PaToH [19] are used for GP and cnHP models, respectively. 555

In the proposed model, we use PaToH as the HP tool in each bipartitioning step. Experiments are conducted with different scaling factors $\alpha = 1, 2, 5$ and 10 for lcn-net cost assignment. We set the maximum allowable imbalance ratio in each bipartitioning as $\epsilon = 0.05$. As both METIS and PaToH involve randomized algorithms in the coarsening phase, five partitioning runs are performed for each instance with different seeds and the averages are reported. We conduct experiments for K = 8, 16, 32, 64, 128and 256 parts (processors).

Table 2 shows the results of the comparison experiments in terms of the communication volume and the reduced system size metrics for dmpGS utilizing the partitions generated by GP, cnHP and the proposed model. For each test instance, these metrics are normalized with respect to the number of rows and the average for all matrices are given for each K. Here and hereafter, all averages are given as geometric means. As seen in Table 2, cnHP achieves considerably low communication volume and reduced system size than GP as expected. The average improvement of cnHP over

				proposed HP model (Sec. 4.1)					
	K	GP	cnHP	$\alpha = 1$	$\alpha = 2$	$\alpha = 5$	$\alpha = 10$		
	8	0.158	0.132	0.140	0.139	0.139	0.145		
ON N	16	0.253	0.217	0.223	0.224	0.224	0.232		
	32	0.380	0.329	0.332	0.332	0.337	0.347		
ц	64	0.547	0.477	0.479	0.479	0.491	0.505		
Con	128	0.767	0.681	0.679	0.680	0.697	0.719		
	256	1.062	0.955	0.948	0.953	0.977	1.012		
ze	8	0.048	0.041	0.033	0.032	0.029	0.029		
SI.	16	0.075	0.066	0.051	0.049	0.045	0.045		
sys.	32	0.109	0.094	0.074	0.070	0.065	0.064		
	64	0.149	0.129	0.102	0.097	0.090	0.088		
ų.	128	0.197	0.174	0.136	0.129	0.119	0.116		
Бe	256	0.252	0.227	0.177	0.168	0.154	0.149		

Table 2: Averages of total communication volume and the reduced system size in dmpGS, both normalized with respect to the number of rows.

GP is approximately 10% for both metrics on K = 256. In fact, cnHP is equivalent 570 to the proposed HP model for $\alpha = 0$. As seen in the table, there is a trade-off 571between the reduced system size and the communication volume for varying values of 572 α for the proposed HP model. Yet the rate of increase in the communication volume 573is observed to be larger than the rate of decrease in the reduced system size with 574increasing α . For example for K = 64, compared to the cnHP model, the proposed 575model slightly increases the communication volume by 0.4%, 0.5%, 2.9% and 5.9%576 whereas it significantly decreases the reduced system size by 21.5%, 25.2%, 30.7% and 577 32.0% for $\alpha = 1, 2, 5$ and 10, respectively. Here, $\alpha = 2$ seems to be a balanced choice 578 since it significantly decreases the reduced system size while it slightly increases the communication volume. This is reflected in the parallel scalability of the proposed 580 algorithm as will be shown in subsection 5.3, thus we set $\alpha = 2$ in the upcoming results. 581 In Figure 7, we provide the performance profiles comparing GP, cnHP and the 582

583 proposed model in terms of the reduced system size. We present the performance profiles only for K = 16,64 and 256 due to lack of space. A performance profile [27] 584shows the comparison of different models relative to the best performing one for each 585 data instance. On a profile, a point (x, y) means that the respective model is within x 586 factor of the best result for a fraction y of the instances. For example, the point (1.20, 587 588 (0.60) on the curve of cnHP means that cnHP yields 20% more reduced system size than the smallest reduced system size achieved for 60% of the dataset. Therefore, the 589model closest to the top left corner is interpreted as the model with best performance. 590As seen in Figure 7, the proposed model outperforms the baseline algorithms in terms of the reduced system size in the majority of the test instances. As K increases,

⁵⁹³ the performance gap between GP and cnHP decreases, whereas the performance gap



Fig. 7: Performance profiles comparing GP, cnHP and the proposed HP model.

kind	ind K = 8		K =	= 16	K =	K = 32		K = 64		K = 128		K = 256	
ID	height	nnz	height	nnz	height	nnz	height	nnz	height	nnz	height	nnz	
1	1,470.1	518.5	554.1	233.6	263.6	143.1	115.9	67.6	65.9	41.1	37.2	25.4	
2	71.9	125.2	63.1	100.6	30.5	61.6	15.8	35.0	8.8	18.7	5.5	11.6	
3	1,219.2	331.4	321.2	271.5	296.8	197.2	167.8	152.7	88.2	82.9	46.1	47.8	
4	27.7	3.8	16.8	5.3	9.9	5.7	8.8	6.2	6.5	4.5	4.6	3.1	
5	260.0	10.0	142.6	9.1	90.3	7.6	63.5	6.3	37.6	4.7	24.5	4.0	
6	600.0	123.2	298.3	101.5	148.6	61.3	73.5	37.0	38.7	22.3	22.0	13.6	
7	131.8	10.1	67.3	7.4	36.7	5.7	22.8	4.7	14.0	3.8	8.5	3.2	
8	513.5	97.3	260.4	59.4	92.2	30.9	48.9	18.5	23.8	11.3	17.0	8.2	
9	1,547.9	1,101.3	1,010.9	1,221.2	556.7	641.8	248.6	315.4	102.0	141.6	50.7	70.0	
10	29.0	4.8	32.9	10.9	15.0	5.6	12.8	5.2	8.6	3.6	6.0	2.7	
11	375.3	619.3	213.3	213.3	112.8	136.6	68.8	89.0	43.2	54.1	25.6	32.0	
12	241.2	170.7	121.4	109.1	59.9	66.8	31.8	41.2	18.1	25.2	10.6	14.7	
13	18.2	2.7	17.7	3.1	12.7	2.7	13.4	2.6	10.1	2.6	8.2	2.8	
14	217.7	231.1	116.5	149.6	59.2	94.0	33.0	54.4	19.1	30.6	12.2	17.7	
15	610.4	228.9	277.9	164.6	122.0	91.6	61.8	50.9	31.5	28.3	18.0	16.4	
16	15.8	62.2	8.5	31.5	6.0	18.9	4.4	11.9	3.4	7.8	2.7	5.3	
All	238.1	57.2	127.1	43.0	65.0	27.8	39.0	18.7	22.7	12.1	14.3	8.3	

Table 3: Total height and nonzero count averages in the off-diagonal blocks of S.

*The values are the ratios of the results attained by the baseline over the proposed in-block reordering.

between the proposed model and both of the baseline models increases significantly. The proposed model yields the best performance for 69%, 71%, 75%, 82%, 85% and 86% of the dataset for K=8, 16, 32, 64, 128 and 256, respectively.

The proposed HP model yields very sparse off-diagonal blocks. The number of nonzeros in any lower off-diagonal block R_k is at most 0.51%, 0.44%, 0.35%, 0.26%, 0.19%, and 0.13% of the total nonzero count of A for K=8, 16, 32, 64, 128, and 256 parts on the average, respectively. As the HP model maintains balance on the nonzero counts of the whole row blocks, these low nonzero counts in off-diagonal blocks do not disturb the computational load balance among processors considerably.

5.2. In-Block Reordering Quality. To our knowledge, no in-block reordering method has been proposed or tested for stSpike in the literature. Therefore, we compare the improvement gained by applying the proposed in-block ordering method against a baseline algorithm which does not apply an in-block reordering. In this comparison, both the proposed and the baseline reordering methods utilize the partitions obtained by the HP model (Section 4.1). Two quality metrics used in this comparison are total height and nonzero count in the off-diagonal blocks of \hat{S} .

Table 3 shows the ratios of these quality metrics of the in-block reorderings generated by the baseline to those of the proposed method. For each K value, the results are given as averages grouped by different matrix kinds, and the last row shows the average of all instances in the dataset.

As seen in Table 3, the proposed reordering method achieves significant improvement in terms of both quality metrics against the baseline reordering. For example for K = 64, on overall average, the proposed method achieves $39 \times$ and $18.7 \times$ improvement against the baseline ordering in terms of height and nonzero counts, respectively. The improvement rate attained in height does not always directly reflect to the improvement rate in the nonzero counts since height is an upper bound for fill-in and the fill-in also depends on the sparsity of the diagonal blocks.

Although the improvement of the proposed reordering against the baseline ordering tends to degrade with increasing K, this is expected since there are fewer rows per block and there is less room for improvement. For example on overall average, the proposed in-block reordering method achieves $57.2 \times$, $43.0 \times$, $27.8 \times$, $18.7 \times$, $12.1 \times$ and $8.3 \times$ decrease in the nonzero count for K = 8, 16, 32, 64, 128 and 256, respectively. The proposed partitioning and reordering model yields very small reduced systems

Matrix	Kind	Sym	Size	Nnz	Relative	$mtGS^*$
11111111	ID	Sym	0120	1112	Residual*	time (s)
msdoor	1	\checkmark	415,863	19,173,163	1.9×10^{-4}	23.1
af_shell1	1	\checkmark	504,855	$17,\!562,\!051$	8.2×10^{-4}	23.4
af_1_k101	1	\checkmark	503,625	17,550,675	1.1×10^{-4}	23.4
CoupCons3D	1		416,800	$17,\!277,\!420$	4.0×10^{-9}	21.8
Freescale1	2		$3,\!428,\!755$	17,052,626	3.0×10^{-4}	72.7
circuit5M_dc	2		3,523,317	14,865,409	1.9×10^{-12}	72.7
CurlCurl_3	9	\checkmark	1,219,574	$13,\!544,\!618$	2.8×10^{-4}	35.4
memchip	2		2,707,524	13,343,948	5.4×10^{-5}	57.5
BenElechi1	6	\checkmark	245,874	$13,\!150,\!496$	6.5×10^{-5}	15.2
pwtk	1	\checkmark	217,918	11,524,432	1.5×10^{-4}	13.6
bmw3_2	1	\checkmark	227,362	11,288,630	1.9×10^{-4}	13.6
bmwcra_1	1	\checkmark	148,770	10,641,602	6.0×10^{-4}	11.9

Table 4: The properties of matrices to run dmpGS.

*Relative residual and runtime results of mtGS on 40 cores for 500 iterations.

whose nonzero counts are significantly low relative to the original system. The average ratios of the nonzero count of the reduced system over the nonzero count of the original coefficient matrix, i.e. $nnz(\hat{S})/nnz(A)$, are 0.05%, 0.12%, 0.26%, 0.49%, 0.87%, and 1.48% for K=8, 16, 32, 64, 128, and 256 parts, respectively. These low nonzero counts of the reduced systems verify the effectiveness of the proposed partitioning and reordering model in terms of alleviating the sequential computational overhead of dmpGS.

5.3. Parallel Scalability. Parallel experiments are performed on the Sariyer
cluster of UHEM [66] using up to 320 cores over 8 distributed nodes, each containing
40 cores (two Intel Xeon Gold 6148 CPUs) and 192GB memory. The nodes are
connected by an InfiniBand EDR 100 Gbps network.

We implement an MPI+OpenMP hybrid parallel dmpGS to demonstrate the ef-638 fectiveness of using stSpike and the proposed model. Throughout this section, the 639 proposed model refers to the proposed partitioning and in-block reordering model 640 641 (Section 4) applied to dmpGS. The number of MPI processes is the same as the number of parts (K) in a partition. For dmpGS, we experimented with different con-642 figurations of number of processes and threads. We found that the best configuration 643 is 8 processes per node and 5 threads per process. Therefore, we conduct parallel 644 experiments for dmpGS using 1, 2, 4 and 8 nodes corresponding to 40, 80, 160 and 645 320 cores and K=8, 16, 32 and 64 parts (processes), respectively. 646

To the best of our knowledge, there is no publicly available true distributedmemory parallel GS implementation. For comparing the performance of dmpGS, we also implemented a multi-threaded GS (mtGS) by using the multithreaded sparse triangular system solver (mkl_sparse_d_trsm) and sparse matrix vector multiplicator (mkl_sparse_d_mv) of Intel MKL [37]. As a baseline, we obtain the results of mtGS on 40 threads/cores (1 node) by using the GP reordering since it is shown in [22] that the triangular solution with MKL benefits most from GP.

We tested the parallel scalability of dmpGS for a subset of the dataset since 654 we have limited core hours on the HPC platform. From the dataset, we considered 655 the matrices with at least 100,000 rows and 10,000,000 nonzeros, for which GS con-656 verges with a relative residual of less than 10^{-3} in 500 iterations with initial guess 657 $x = [0, \ldots, 0]^T$ and right-hand side vector $f = [1/m, 2/m, \ldots, 1]^T$. Then we select only 658 those instances with different sparsity structures from each matrix group. There were 659 exactly 12 such matrices in our dataset satisfying these criteria. The properties of 660 those matrices are shown in Table 4, sorted in decreasing order of nonzero counts. 661 662 The sixth and the last column respectively show the relative residual and runtime of

	number of		CD UD		proposed model			
K	nodes	cores	GP	cnHP	$\alpha = 1$	$\alpha = 2$	$\alpha = 5$	$\alpha = 10$
8	1	40	9.87	8.71	14.85	14.71	14.77	14.51
16	2	80	14.65	13.58	29.07	28.51	28.47	28.25
32	4	160	17.41	16.11	47.28	47.86	47.24	45.89
64	8	320	15.79	17.60	54.96	55.54	50.21	50.65

Table 5: Average speedup obtained by dmpGS over mtGS on 40 cores.

*The best speedup value obtained for each K is shown in bold.

663 mtGS after 500 iterations.

Table 5 shows the average speedup values obtained by dmpGS with GP, cnHP and the proposed model over mtGS. We run dmpGS with the proposed model for $\alpha = 1, 2, 5$ and 10 to observe the effect of scaling factor (α) on the parallel performance. As seen in the table, the proposed model achieves significantly higher speedup for dmpGS over the baseline models for all α . The speedup performance gap between the proposed and baseline models increase with increasing K, thus confirming the effectiveness of the proposed model.

We also provide Figure 8 which depicts the performance profiles for comparing the dmpGS runtime using the proposed model for varying α and K values. We choose $\alpha=2$ for better scalability of dmpGS since it yields the best performance for larger part counts (K=32 and 64) as seen in both Table 5 and Figure 8. As seen in Table 5, the proposed model with $\alpha=2$ yields average of $1.5 \times$, $1.9 \times$, $2.7 \times$ and $3.2 \times$ higher speedup relative to the best of the baseline models for K=8, 16, 32 and 64, respectively.

Figure 9 shows the results of the strong scaling experiments as speedup curves of dmpGS with GP, cnHP and the proposed model. The proposed model significantly enhances the scalability of dmpGS so that dmpGS scales up to 320 cores on all instances. As seen in the figure, the proposed model outperforms GP and cnHP models for all of the test instances, significantly so in 9 out of 12. In Figure 9 for memchip, dmpGS using the proposed model achieves up to 122.2 speedup on 320 cores over mtGS on 40 cores.

6. Conclusion. We proposed and implemented an stSpike-based distributed-684 memory parallel GS (dmpGS) algorithm. For improving the scalability of dmpGS, 685 we propose a hypergraph partitioning (HP) based partitioning model and an in-block 686 row reordering method. Extensive experiments show that the proposed HP model 687 688 significantly decreases the reduced system size with respect to the baseline models while attaining comparable communication volume. The proposed in-block reordering 689 method leads to a substantial decrease in the computational cost of both forming and 690 solving the reduced system. Parallel experiments up to 320 cores demonstrate that 691 using the proposed reordering model significantly improves the scalability of dmpGS. 692



Fig. 8: Performance profiles in terms of the dmpGS runtime using the proposed model.



Fig. 9: Speedup curves of dmpGS with GP, cnHP and the proposed model (for K=8, 16, 32 and 64) relative to mtGS on 1 node (40 cores).

As a future work, we will consider the parallel solution of the reduced system 693 to further alleviate the sequential bottleneck. We will also consider an in-block row 694 reordering which takes the nonzeros of the diagonal blocks into account for further 695 reducing the nonzero count in the reduced system. Finally, the future work will 696 697 include extending the dmpGS algorithm for multiple right-hand-side vectors as it is very common in modern applications. Using multiple right-hand-side vectors is 698 expected to further enhance the performance of dmpGS since it enables using higher 699 level BLAS subroutines compared to the single right-hand-side case. Moreover, the 700 parallel solution time per right-hand-side vector will further decrease since the parallel 701 702 factorization is done only once.

T. TORUN, F. S. TORUN, M. MANGUOGLU, AND C. AYKANAT

703		REFERENCES
704	[+1	CATTER F. Kuller en euro C. Automatic Alexandre en difference del forma filosofici
704 705	[1]	5. ACER, E. KAYAASLAN, AND C. AYKANAT, A nypergraph partitioning model for profile mini- mization SIAM Journal on Scientific Computing 41 (2010) pp. 483–4108
705	[2]	I. M. ADAMS AND H. F. LORDAN. Is SOR color-blind? SIAM Journal on Scientific and Statis-
707	[2]	tical Computing, 7 (1986), pp. 490–506.
708	[3]	M. ADAMS, M. BREZINA, J. HU, AND R. TUMINARO, Parallel multigrid smoothing: polynomial
709		versus Gauss-Seidel, Journal of Computational Physics, 188 (2003), pp. 593-610.
710	[4]	M. F. ADAMS, A distributed memory unstructured Gauss-Seidel algorithm for multigrid
711		smoothers, in SC'01: Proceedings of the 2001 ACM/IEEE Conference on Supercomputing,
712	[]	2001.
713	[5]	P. R. AMESTOY, I. S. DUFF, JY. L'EXCELLENT, AND J. KOSTER, A fully asynchronous multi-
714 715		frontal solver using distributed dynamic scheduling, SIAM Journal on Matrix Analysis and
716	[6]	P AMODIO AND F MAZZIA A narallel Gauss-Seidel method for block tridiagonal linear sus-
717	[0]	tems. SIAM Journal on Scientific Computing, 16 (1995), pp. 1451–1461.
718	[7]	R. BAGNARA, A unified proof for the convergence of Jacobi and Gauss-Seidel methods, SIAM
719		review, 37 (1995), pp. 93–97.
720	[8]	G. BALLARD, A. DRUINSKY, N. KNIGHT, AND O. SCHWARTZ, Hypergraph partitioning for sparse
721		$matrix-matrix\ multiplication,\ ACM\ Transactions\ on\ Parallel\ Computing\ (TOPC),\ 3\ (2016),$
722	[0]	pp. 1–34.
(23 794	[9]	R. BARRETT, M. BERRY, I. F. CHAN, J. DEMMEL, J. DONATO, J. DONGARRA, V. EIJKHOUT, B. DOZO, C. BOMME, AND H. VAN DEP VOPET. Templates for the solution of linear systems:
725		huilding blocks for iterative methods SIAM 1994
726	[10]	M. BENZI, Preconditioning techniques for large linear systems: a survey, Journal of Computa-
727	r - 1	tional Physics, 182 (2002), pp. 418–477.
728	[11]	M. BENZI AND M. TUMA, A sparse approximate inverse preconditioner for nonsymmetric linear
729		systems, SIAM Journal on Scientific Computing, 19 (1998), pp. 968–994.
730	[12]	M. BERNASCHI, P. D'AMBRA, AND D. PASQUINI, AMG based on compatible weighted matching
731	[10]	for GPUs, Parallel Computing, 92 (2020).
(32 733	[13]	M. W. BERRY, Large-scale sparse singular value computations, The International Journal of Supercomputing Applications 6 (1002) pp. 13-40
734	[14]	B H BISSELING Parallel Scientific Computation: A Structured Approach Using BSP Oxford
735	[1 1]	University Press, USA, 2020.
736	[15]	E. S. BOLUKBASI AND M. MANGUOGLU, A multithreaded recursive and nonrecursive parallel
737		sparse direct solver, in Advances in Computational Fluid-Structure Interaction and Flow
738		Simulation, Springer, 2016, pp. 283–292.
739	[16]	J. BOYLE, M. MIHAJLOVIC, AND J. SCOTT, HSL M120: an efficient AMG preconditioner, tech.
740 741	[17]	report, Citeseer, 2007. B. BYLINA AND I. BYLINA Analysis and comparison of mondaring for two factorization methods
741	[17]	(LU and WZ) for snarse matrices in International Conference on Computational Science
743		Springer, 2008, pp. 983–992.
744	[18]	U. V. CATALYUREK AND C. AYKANAT, Hypergraph-partitioning-based decomposition for parallel
745		sparse-matrix vector multiplication, IEEE Transactions on parallel and distributed systems,
746		10 (1999), pp. 673–693.
747	[19]	U. V. ÇATALYÜREK AND C. AYKANAT, PaToH (Partitioning Tool for Hypergraphs), in Ency-
748	[20]	clopedia of Parallel Computing, Springer, 2011, pp. 1479–1487.
749	[20]	U. V. GATALYUREK, U. AYKANAT, AND B. UÇAR, On two-aimensional sparse matrix partition-
751		ng. Models, methods, and a recipe, SIAM Journal on Scientific Computing, 32 (2010),
752	[21]	SC. CHEN, D. J. KUCK, AND A. H. SAMEH, Practical parallel band triangular system solvers,
753		ACM Transactions on Mathematical Software (TOMS), 4 (1978), pp. 270–277.
754	[22]	I. ÇUĞU AND M. MANGUOĞLU, A parallel multithreaded sparse triangular linear system solver,
755		Computers & Mathematics with Applications, 80 (2020), pp. 371–385.
756	[23]	T. A. DAVIS, Direct methods for sparse linear systems, SIAM, 2006, pp. 83–95.
757 759	[24]	T. A. DAVIS AND Y. HU, The University of Florida sparse matrix collection, ACM Transactions
759	[25]	K D DEVINE E G BOMAN R T HEAPHY R H RISSELING AND U V CATALVUDER
760	[40]	Parallel hyperaraph partitioning for scientific computing, in Proceedings 20th IEEE Inter-
761		national Parallel & Distributed Processing Symposium, IEEE, 2006, pp. 10–pp.
762	[26]	J. DÍAZ, J. PETIT, AND M. SERNA, A survey of graph layout problems, ACM Computing Surveys
763		(CSUR), 34 (2002), pp. 313–356.

22

- [27] E. D. DOLAN AND J. J. MORÉ, Benchmarking optimization software with performance profiles,
 Mathematical programming, 91 (2002), pp. 201–213.
- [28] J. J. DONGARRA AND A. H. SAMEH, On some parallel banded system solvers, Parallel Computing, 1 (1984), pp. 223–235.
- [29] I. S. DUFF AND J. KOSTER, On algorithms for permuting large entries to the diagonal of a sparse matrix, SIAM Journal on Matrix Analysis and Applications, 22 (2001), pp. 973–996.
- [30] H. C. ELMAN, D. J. SILVESTER, AND A. J. WATHEN, Finite elements and fast iterative solvers:
 with applications in incompressible fluid dynamics, Numerical Mathematics and Scientific
 Computation, 2014.
- [31] G. C. Fox, Solving problems on concurrent processors, Old Tappan, NJ; Prentice Hall Inc.,
 1988.
- [32] M. GEE, C. SIEFERT, J. HU, R. TUMINARO, AND M. SALA, *ML 5.0 smoothed aggregation user's guide*, Tech. Report SAND2006-2649, Sandia National Laboratories, 2006.
- [33] G. GOLUB AND J. M. ORTEGA, Scientific computing: an introduction with parallel computing,
 Academic Press Professional, Inc., 1993.
- [34] G. H. GOLUB AND C. F. VAN LOAN, Matrix computations, vol. 3, JHU press, 2013, pp. 606–616.
- [35] L. GRIGORI, S. MOUFAWAD, AND F. NATAF, Enlarged Krylov subspace conjugate gradient methods for reducing communication, SIAM Journal on Matrix Analysis and Applications, 37 (2016), pp. 744–773.
- [36] V. E. HENSON AND U. M. YANG, BoomerAMG: A parallel algebraic multigrid solver and preconditioner, Applied Numerical Mathematics, 41 (2002), pp. 155–177.
- [37] INTEL, Intel Math Kernel Library (MKL). https://software.intel.com/en-us/mkl, 2019.
- [38] K. S. KANG, Scalable implementation of the parallel multigrid method on massively parallel computers, Computers & Mathematics with Applications, 70 (2015), pp. 2701–2708.
- [39] G. KARYPIS, R. AGGARWAL, V. KUMAR, AND S. SHEKHAR, Multilevel hypergraph partitioning:
 Applications in VLSI domain, IEEE Transactions on Very Large Scale Integration (VLSI)
 Systems, 7 (1999), pp. 69–79.
- [40] G. KARYPIS AND V. KUMAR, METIS: Unstructured Graph Partitioning and Sparse Matrix
 Ordering System, Version 5.1. http://www.cs.umn.edu/~metis, 2013.
- [41] M. KAWAI, T. IWASHITA, H. NAKASHIMA, AND O. MARQUES, Parallel smoother based on block red-black ordering for multigrid Poisson solver, in International Conference on High Performance Computing for Computational Science, Springer, 2012, pp. 292–299.
- [42] D. P. KOESTER, S. RANKA, AND G. C. FOX, A parallel Gauss-Seidel algorithm for sparse power
 system matrices, in Supercomputing'94: Proceedings of the 1994 ACM/IEEE Conference
 on Supercomputing, IEEE, 1994, pp. 184–193.
- [43] X. S. LI AND J. W. DEMMEL, SuperLU_DIST: A scalable distributed-memory sparse direct solver for unsymmetric linear systems, ACM Transactions on Mathematical Software (TOMS), 29 (2003), pp. 110–140.
- [44] Y. LIN AND J. YUAN, Profile minimization problem for matrices and graphs, Acta Mathemat icae Applicatae Sinica, 10 (1994), pp. 107–112.
- [45] W. LIU, A. LI, J. HOGG, I. S. DUFF, AND B. VINTER, A synchronization-free algorithm for parallel sparse triangular solves, in European Conference on Parallel Processing, Springer, 2016, pp. 617–630.
- [46] L. O. MAFTEIU-SCAI, The bandwidths of a matrix. a survey of algorithms, Annals of West
 University of Timisoara-Mathematics, 52 (2014), pp. 183–223.
- [47] M. MANGUOGLU, A domain-decomposing parallel sparse linear system solver, Journal of Com putational and Applied Mathematics, 236 (2011), pp. 319–325.
- [48] M. MANGUOGLU, Parallel solution of sparse linear systems, in High-Performance Scientific
 Computing, Springer, 2012, pp. 171–184.
- [49] M. MANGUOGLU, A. H. SAMEH, AND O. SCHENK, PSPIKE: A parallel hybrid sparse linear
 system solver, in European Conference on Parallel Processing, Springer, 2009, pp. 797–
 808.
- [50] C. C. K. MIKKELSEN AND M. MANGUOGLU, Analysis of the truncated SPIKE algorithm, SIAM
 Journal on Matrix Analysis and Applications, 30 (2009), pp. 1500–1519.
- [51] Y. NOTAY, An aggregation-based algebraic multigrid method, Electronic transactions on numer ical analysis, 37 (2010), pp. 123–146.
- [52] D. P. O'LEARY, Ordering schemes for parallel processing of certain mesh problems, SIAM
 Journal on Scientific and Statistical Computing, 5 (1984), pp. 620–632.
- [53] W. PAZNER, Efficient low-order refined preconditioners for high-order matrix-free continuous and discontinuous Galerkin methods, SIAM Journal on Scientific Computing, 42 (2020), pp. A3055–A3083.
- 825 [54] E. POLIZZI AND A. SAMEH, A parallel hybrid banded system solver: the SPIKE algorithm,

T. TORUN, F. S. TORUN, M. MANGUOGLU, AND C. AYKANAT

826 Parallel computing, 32 (2006), pp. 177–194.

24

- [55] E. POLIZZI AND A. SAMEH, SPIKE: A parallel environment for solving banded linear systems,
 Computers & Fluids, 36 (2007), pp. 113–120. Challenges and Advances in Flow Simulation
 and Modeling.
- [56] E. POLIZZI AND A. H. SAMEH, A parallel hybrid banded system solver: the SPIKE algorithm,
 Parallel Computing, 32 (2006), pp. 177–194. Parallel Matrix Algorithms and Applications
 (PMAA'04).
- [57] J. W. RUGE AND K. STÜBEN, Algebraic multigrid, in Multigrid methods, SIAM, 1987, pp. 73–
 130.
- [58] Y. SAAD, *ILUT: A dual threshold incomplete LU factorization*, Numerical linear algebra with
 applications, 1 (1994), pp. 387–402.
- 837 [59] Y. SAAD, Iterative methods for sparse linear systems, SIAM, 2003.
- [60] A. H. SAMEH AND R. P. BRENT, Solving triangular systems on a parallel computer, SIAM
 Journal on Numerical Analysis, 14 (1977), pp. 1101–1113.
- [61] O. SCHENK, M. MANGUOGLU, A. SAMEH, M. CHRISTEN, AND M. SATHE, Parallel scalable PDEconstrained optimization: antenna identification in hyperthermia cancer treatment planning, Computer Science-Research and Development, 23 (2009), pp. 177–183.
- [62] Y. SHANG, A distributed memory parallel Gauss-Seidel algorithm for linear algebraic systems,
 Computers & Mathematics with Applications, 57 (2009), pp. 1369–1376.
- [63] B. S. SPRING, E. POLIZZI, AND A. H. SAMEH, A feature-complete SPIKE dense banded solver,
 ACM Transactions on Mathematical Software (TOMS), 46 (2020), pp. 1–35.
- [64] R. TAVAKOLI AND P. DAVAMI, A new parallel Gauss-Seidel method based on alternating group
 explicit method and domain decomposition method, Applied mathematics and computation,
 188 (2007), pp. 713–719.
- [65] B. UÇAR AND C. AYKANAT, Revisiting hypergraph models for sparse matrix partitioning, SIAM
 review, 49 (2007), pp. 595–603.
- [66] UHEM, National Center for High Performance Computing. http://www.uhem.itu.edu.tr, 2021.
- [67] P. VAN, M. BREZINA, J. MANDEL, ET AL., Convergence of algebraic multigrid based on smoothed
 aggregation, Numerische Mathematik, 88 (2001), pp. 559–579.
- [68] B. VASTENHOUW AND R. H. BISSELING, A two-dimensional data distribution method for parallel sparse matrix-vector multiplication, SIAM review, 47 (2005), pp. 67–95.
- [69] I. E. VENETIS, A. KOURIS, A. SOBCZYK, E. GALLOPOULOS, AND A. H. SAMEH, A direct tridiagonal solver based on Givens rotations for GPU architectures, Parallel Computing, 49 (2015), pp. 101–116.
- [70] K. WATANABE, H. IGARASHI, AND T. HONMA, Comparison of geometric and algebraic multigrid methods in edge-based finite-element analysis, IEEE transactions on magnetics, 41 (2005), pp. 1672–1675.
- [71] R. WEBSTER, An algebraic multigrid solver for Navier-Stokes problems, International Journal
 for Numerical Methods in Fluids, 18 (1994), pp. 761–780.
- [72] U. M. YANG, Parallel algebraic multigrid methods-high performance preconditioners, in Nu merical solution of partial differential equations on parallel computers, Springer, 2006,
 pp. 209–236.
- [73] J. ZHANG, Acceleration of five-point red-black Gauss-Seidel in multigrid for Poisson equation,
 Applied Mathematics and Computation, 80 (1996), pp. 73–93.