

**Abstract.** This is an example SIAM L<sup>A</sup>T<sub>E</sub>X article. This can be used as a template for new articles. Abstracts must be able to stand alone and so cannot contain citations to the paper's references, equations, etc. An abstract must consist of a single paragraph and be concise. Because of online formatting, abstracts must appear as plain as possible. Any equations should be inline.

**Key words.** example, L<sup>A</sup>T<sub>E</sub>X

**AMS subject classifications.** 68Q25, 68R10, 68U05

**1. Introduction.** The introduction introduces the context and summarizes the manuscript. It is important to clearly state the contributions of this piece of work. The next two paragraphs are text filler, generated by the `lipsum` package.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

The paper is organized as follows. Our main results are in ??, our new algorithm is in ??, experimental results are in ??, and the conclusions follow in ??.

**2. Related Work.** One of the earliest work on parallelizing kernels having loop-carried dependencies is the red-black Gauss-Seidel scheme [6]. Later Kamath and Sameh introduced a two-block partitioning scheme for parallelizing Kaczmarz method on tridiagonal structures [13]. A general study on the convergence of these block methods were done by Elfving in 1980 [5].

The advent of processors having more parallelism and the need to consider more unstructured matrices have made graph-based approach an important tool for parallelizing such kernels. Multicoloring is one of the most popular approach used in this field [12], but is sometimes not efficient on modern cache-based processors. There have been several researches going on to increase the efficiency of multi-coloring and improving the heuristics, an overview of the methods can be found in [15]. One of the most successful method in this regard is the algebraic block multi-coloring [11] proposed by Iwashita et al. in 2012.

Another line of research focuses on parallelizing dependent kernels while maintaining the same convergence behavior of sequential execution. One of the earliest known

---

\*Submitted to the editors DATE.

**Funding:** This work was funded by the Fog Research Institute under contract no. FRI-454.

works in this category is the hyperplane method [21]. Extensions to this approach can be seen in [17] where a hybrid approach between multi-coloring and hyperplane method is used. Most recent work in this direction can be attributed to Park et al. on his work with level-scheduling for triangular solvers [18].

Most of the above mentioned method have been tested only for their applicability to parallelize distance-1 dependent kernels and some of them are not capable to deal with dependencies like distance-2. The research on parallelizing distance-1 dependent kernels has been strongly accelerated after the introduction of HPCG benchmark [4]. When it comes to distance-2 kernels popular methods seen in the literature are locking based methods, thread private local vectors [8] for kernels like symmetric sparse matrix vector or with the usage of specially tailored sparse matrix data formats like compressed sparse blocks (CSB) [3] or recursive sparse blocks (RSB) [16].

**3. Contribution.** The paper focuses on developing an alternative method to parallelize kernels having loop-carried dependencies. The method introduced here is applicable for solving general distance- $k$  dependencies, similar to multi-coloring methods. Currently we focus only on undirected graph i.e., matrices with symmetric pattern. The main motivation of the approach is to achieve good hardware performance on modern hardware architecture, by generating sufficient parallelism while preserving good data locality. The method needs no specialized data format, and works basically on simple sparse matrix format like compressed row storage (CRS).

Most of the above approaches explained above in section 2 suffer from performance penalties in one way or the other, for example multi-coloring degrades the data locality, although this can be improved considerably using algebraic block multi-coloring, still for moderately large matrices or with the increase in  $k$  of distance- $k$  dependency the method shows deterioration in performance. Similar drawbacks exists for other methods which will be discussed in detail within this paper.

In this work we provide a detailed performance analysis of the method and comparison between different existing methods chosen from representative classes. The comparisons are done both for exact kernels like symmetric sparse matrix vector (SymmSpMV) having distance-2 dependency and iterative solvers like Gauss-Seidel (GS) and Kaczmarz (KACZ) schemes having distance-1 and distance-2 dependencies respectively. For iterative schemes we further provide comparison between convergence of different methods. The comparisons are done on different hardware architectures ranging from Intel's Ivy-Bridge series to modern Sky-Lake architecture and the AMD Epyc architecture. The comparisons shows the superiority of our method compared to others and the applicability of our method on wide-variety of heterogeneous systems. As far to our knowledge this is the first paper which demonstrates such high efficiency of distance-2 dependent kernels using simple and common CRS matrix storage format on such broad scale of matrices.

The paper is limited to socket level, and we use only thread level parallelization. Multi-node parallelization is left for future work. However it should be noted that for iterative kernels like KACZ and GS node-level performance is far more important because commonly such solvers are applied only locally and a different approaches are used for parallelizing between nodes [4, 10].

As a final application run we demonstrate the parallelization of an eigen-value solver called FEAST [19], where we use an iterative inner linear solver based on Kaczmarz method. The result presented is the first to achieve such high performance on node level for an iterative solver and is superior to the previous results published [7].

#### 4. Test bed, matrices and kernels.

**4.1. Test bed.** The tests are conducted on three different multi-core architectures. Two of them being Intel’s Ivy-Bridge and modern Sky-Lake architecture, the choice of these architectures enables a good comparison of our method on both Low Core Count (LCC) and High Core Count (HCC) architectures, which would play an important role in the convergence of the method we develop. As a third choice we select AMD’s recent Epyc architecture, which is competitive to Intel Sky-Lake architecture. This choice enables us to study the effect of our method on chips based on completely different microarchitecture, enabling us to demonstrate the applicability of our method on wide range of architectures. All the tests are conducted on a single socket of these architectures.

- Intel Ivy-Bridge architecture belongs to class of classic Intel’s cache-based architecture, which has three inclusive cache hierarchies. All the cache are scalable and the LLC (L3) being shared among all the cores on one socket. The processor is capable of delivering one full four wide SIMD add, multiply and load in one cycle.
- Intel Sky-Lake architecture belongs to recent generation of Intel family. Contrary to it’s predecessors (like Ivy-Bridge ), the LLC is now changed to a non-inclusive victim cache shared by all the cores on a socket. The architecture comes with support for eight wide SIMD operations (AVX-512). The processor is capable of doing two AVX-512 add, multiply and load operations per cycle.
- AMD Epyc is based on AMD’s Zen microarchitecture. The basic building block of the architecture consists of Core Complex (CCX) consisting of three cores (can extend upto four on high end models) each having it’s own private L1 and L2 cache. The L3 cache is shared between a core complex and is non-inclusive victim cache. A single socket of Epyc consists of eight such core complexes.

The details of architectures along with the measured bandwidths are given in [Table 1](#). The bandwidths are measured using *likwid – bench* suite.

TABLE 1  
Test bed

Model name	Xeon <sup>®</sup> E5-2660	Xeon <sup>®</sup> Gold 6148	Epyc 7451
Microarchitecture	Ivy Bridge	Skylake	Zen
Clock	2.2 GHz	2.4 GHz	2.3 GHz
Physical Cores per socket	10	20	24
L1d Cache	10 × 32 kB	20 × 32 kB	24 × 32 kB
L2 Cache	10 × 256 kB	20 × 1 MB	24 × 512 MB
L3 Cache	25 MB	27.5 MB	8 × 8 MB
L3 type	inclusive	non-inclusive	non-inclusive
Main Memory	32 GB	45 GB	4 × 16 GB
Bandwidth - load only)	47 GB/s	110 GB/s	130 GB/s
Bandwidth - copy	40 GB/s	110 GB/s	114 GB/s
Architecture specific flag	-	-xCORE-AVX512	-

The code was compiled with newest Intel compiler version 18 and the following compiler flags were set `-fno-alias -xHost -O3`. Furthermore all the measurements were done with CPU clock speeds fixed at frequencies indicated in [Table 1](#).

**4.2. Benchmark Matrices.** All the test matrices are taken from SuiteSparse Matrix Collection (former University of Florida Sparse Matrix Collection) [2] and ESSEX project [1]. The selection of the matrices from SuiteSparse Matrix Collection is mainly done by combining the test matrices from two papers [16, 18]. This enables easy comparison of results. Matrices from ESSEX project are some of the matrices that are of interest in the FEAST eigen value solver. Only matrices having undirected graphs are considered due to scope of the paper as mentioned in section 3. Matrices along with some of their parameters are given in Table 2. Matrices that have been marked with an \* symbol indicate they are corner cases and will be discussed in detail.

TABLE 2  
Benchmark matrices

Index	Matrix name	nrows	nnz	bandwidth		
1	audikw_1	943695	77651847	925946		
2	bone010	986703	71666325	13016		
3	channel-500x100x100-b050	4802000	85362744	600299		
4	crankseg_1	52804	10614210	50388	*	
5	delaunay_n24	16777216	100663202	16769102		
6	dielFilterV3real	1102824	89306020	1036475		
7	Emilia_923	923136	41005206	17279		
8	F1	343791	26837113	343754		
9	Fault_639	638802	28614564	19988		
10	Flan_1565	1564794	117406044	20702		
11	G3_circuit	1585478	7660826	947128		
12	Geo_1438	1437960	63156690	26018		
13	gsm_106857	589446	21758924	588744		
14	Hook_1498	1498023	60917445	29036		
15	HPCG-192	7077888	189119224	37057		
16	inline_1	503712	36816342	502403		
17	nlpkkt120	3542400	96845792	1814521		
18	nlpkkt200	16240000	448225632	8240201	*	
19	offshore	259789	4242673	237738	*	
20	parabolic_fem	525825	3674625	525820	*	
21	pwt_k	217918	11634424	189331		
22	Serena	1391349	64531701	81578		
23	ship_003	121728	8086034	3659	*	
24	thermal2	1228045	8580313	1226000	*	
25	Anderson-16.5	2097152	14680064	1198372		
26	Graphene-2048	4194304	16771072	2048		
27	Graphene-4096	16777216	67096576	4096		
28	Spin-26	10400600	145608400	709995	*	

**4.3. Kernels.** To test the performance we choose algorithms that are exact as well as iterative. Also we include kernels from both distance-1 and distance-2 dependency classes. All the kernels shown below are based on CRS matrix storage format.

**4.3.1. SpMV .** Sparse Matrix Vector (SpMV ) is a kernel that do not have any dependencies. It acts as a good reference for other kernels to determine their performance upper bound.

**Algorithm 4.1** SpMV Find  $b : b = Ax$ 


---

```

1: for  $row = 1 : n_{rows}$  do
2:   for  $idx = rowPtr[row] : rowPtr[row + 1]$  do
3:      $b[row] += A[idx] * x[col[idx]]$ 
4:   end for
5: end for

```

---

142 The arithmetic intensity of the kernel  $I_{SpMV}$  is as follows:

143 (4.1) 
$$I_{SpMV} = \frac{2}{8 + 4 + 8 * \alpha + \frac{16}{N_{nzt}}}$$

144 where  $\alpha$  represents the data locality factor and  $N_{nzt}$  non-zeros per row.  $\alpha$  depends  
 145 on the sparsity pattern of the matrix and varies from matrix to matrix. Ideal value  
 146 of  $\alpha$  for sufficiently large matrix is  $\frac{1}{N_{nzt}}$ . More details on factor  $\alpha$  could be found in  
 147 [14].

148 **4.3.2. SpMTV .** Sparse Matrix Transpose Vector (SpMTV ) is a kernel having  
 149 distance-2 dependency.

**Algorithm 4.2** SpMTV Find  $b : b = A'x$ 


---

```

1: for  $row = 1 : n_{rows}$  do
2:   for  $idx = rowPtr[row] : rowPtr[row + 1]$  do
3:      $b[col[idx]] += A[idx] * x[row]$ 
4:   end for
5: end for

```

---

150 In comparison to SpMV operation, the kernel requires an extra scatter operation,  
 151 which causes dependency. The arithmetic intensity of the kernel  $I_{SpMTV}$  is given as:

152 (4.2) 
$$I_{SpMTV} = \frac{2}{8 + 4 + 16 * \alpha + \frac{8}{N_{nzt}}}$$

153 In ideal case data traffic for this kernel should remain close to that of SpMV, if  
 154  $N_{nzt}$  are sufficiently high, and  $\alpha$  factor is small enough.

155 **4.3.3. SymmSpMV .** Symmetric Sparse Matrix Vector (SymmSpMV ) makes  
 156 use of the symmetric property of the matrix to perform the matrix vector multiplica-  
 157 tion.

**Algorithm 4.3** SymmSpMV Find  $b : b = Ax$ , where  $A$  is an upper triangular matrix

---

```

1: for  $row = 1 : n_{rows}$  do
2:    $diag\_idx = rowPtr[row]$ 
3:    $b[row] += A[diag\_idx] * x[row]$ 
4:   for  $idx = rowPtr[row] + 1 : rowPtr[row + 1]$  do
5:      $b[row] += A[idx] * x[col[idx]]$ 
6:      $b[col[idx]] += A[idx] * x[row]$ 
7:   end for
8: end for

```

---

To operate on this kernel we just use the upper triangular part of the sparse matrix. The kernel requires only half the data traffic compared to SpMV but requires the same amount of Flops, leading to almost twice the intensity of SpMV operations.

$$(4.3) \quad I_{SymmSpMV} = \frac{4}{8 + 4 + 32 * \alpha + \frac{16}{N_{nzt}^{symm}}}$$

Note that  $N_{nzt}^{symm}$  is the number of non-zeros per row in upper triangular part of the matrix.

**4.3.4. GS .** Gauss-Seidel (GS ) is a solver having distance-1 dependency. Contrary to the above kernels GS is in-exact meaning it is an iterative method. Alg. 4.4 shows the Gauss-Seidel algorithm where its assumed that the diagonal entries of the matrix is stored as first entry in their corresponding rows.

---

**Algorithm 4.4** GS Solve for  $x : Ax = b$

---

```

1: for row = 1 : rows do
2:   x[row]+ = b[row]
3:   for idx = rowPtr[row] + 1 : rowPtr[row + 1] do
4:     x[row]- = A[idx] * x[col[idx]]
5:   end for
6:   diag = A[rowPtr[row]]
7:   x[row]/ = diag
8: end for

```

---

Regarding the in-core execution the kernel has same properties as of SpMV , but requires an additional divide operation per row of the matrix. If the locality ( $\alpha$  factor) is not disturbed due to pre-processing the kernel requires same data traffic as of SpMV . The arithmetic intensity of GS is the same as that of SpMV , if we neglect the divide operation that occurs once per every row.

$$(4.4) \quad I_{GS} = I_{SPMV}$$

**4.3.5. KACZ.** Kaczmarz (KACZ ) is an iterative solver based on row-projection based methods. The solver has a distance-2 dependency.

---

**Algorithm 4.5** KACZ Solve for  $x : Ax = b$

---

```

1: for row = 1 : rows do
2:   row_norm = 0
3:   scale = b[row]
4:   for idx = rowPtr[row] : rowPtr[row + 1] do
5:     scale- = A[idx] * x[col[idx]]
6:     row_norm+ = A[idx] * A[idx]
7:   end for
8:   scale = scale/row_norm
9:   for idx = rowPtr[row] : rowPtr[row + 1] do
10:    x[col[idx]]+ = scale * A[idx]
11:   end for
12: end for

```

---

In-core has a mixed behavior of both SpMV and SpMTV similar to SymmSpMV. The solver also requires a divide per row of the matrix. In ideal case the data traffic

from memory should remain same as that of SpMTV. But the solver requires thrice the flops compared to SpMTV per non-zero. For brevity of the results we ignore the flops used in *rownorm* computations since, one could also row normalize the sparse matrix before performing the KACZ operation. This leads to an almost two fold higher Arithmetic Intensity compared to SpMTV.

$$(4.5) \quad I_{KACZ} = \frac{4}{8 + 4 + 16 * \alpha + \frac{8}{nnzr}} = 2 * I_{SpMTV}$$

**5. Motivation.** Motivation for developing an alternative method stems from the ESSEX (Equipping Sparse Solvers for Exascale) project [1] where we investigate into solving large eigen-value problems from quantum mechanics field. In this context having a robust iterative solver was inevitable, due to the poor condition number of the matrices that appear in this field. Kaczmarz (KACZ) solver was found to be satisfactory but parallelizing this solver was deemed challenging because of the loop-carried dependencies in the kernel. Previous work on parallelizing the KACZ kernel used multi-coloring (MC) [7] but it was soon found that the kernels do not scale efficiently with this approach.

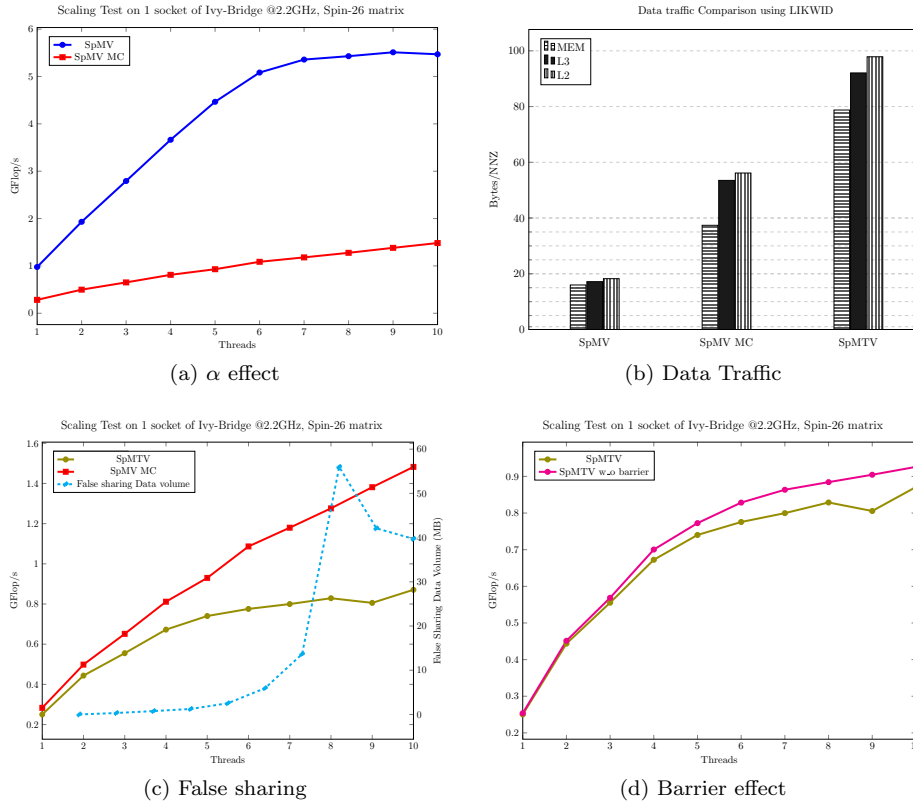


FIG. 1. *Effect of Multicoloring*

In order to get a better understanding of the underlying problem it's convenient to choose simple sparse matrix transpose vector (SpMTV) as a benchmark kernel. The particular choice of this kernel is due to the fact that both KACZ and SpMTV

have similar kind of dependencies, and it's much easier to compare with our reference kernel namely sparse matrix vector (SpMV) which is embarrassingly parallel. The algorithm for SpMTV and SpMV has been listed in Algorithms 4.1 and 4.2

Figure 1a shows the performance of SpMV kernel on original unpermuted matrix and matrix with MC permutation. Here we see the performance of SpMV on multi-colored matrix is five times worse than that of SpMV on unpermuted matrix. One of the major reason for this drop is due to the increase in  $\alpha$  factor seen in the intensity equation (4.2) Since the kernels like SpMV are mainly memory bound increase in  $\alpha$  lowers intensity  $I_{SpMV}$  leading to a drop of performance as predicted by roofline model [23]. This could easily be demonstrated by measuring the data traffic between different memory hierarchies. We do this using the LIKWID tool [22], and the measurements can be seen in Figure 1b. One can see an increase in data-traffic from all the memory hierarchy compared to SpMV on normal unpermuted matrix. This is basically caused by the bad data locality introduced by multi-coloring permutation.

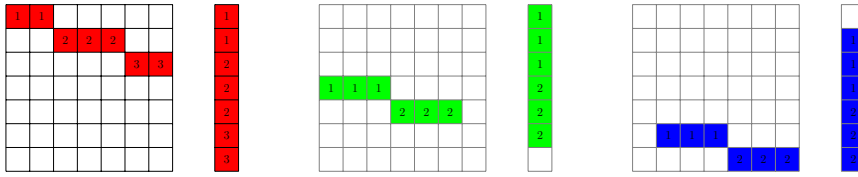


FIG. 2. Illustration of increase in  $\alpha$  by multicoloring, numbers represents thread numbers working on a particular row

Figure 2 shows an illustration of why data traffic increases for a given matrix. If one assumes last level cache (LLC) can only hold less than six elements and obeys perfect LRU policy, as seen in the Figure 2 for each new color we would need to load the data from main memory. As we will see later this  $\alpha$  factor strongly depends on the matrix size and the size of LLC.

As seen in Figure 1b the data traffic further increases for SpMTV due to additional indirect writes (scatter) and this scales up  $\alpha$  factor as seen in the denominator of  $I_{SpMTV}$  (see (4.2)). This further decreases performance of SpMTV compared to SpMV on MC matrix which can be verified in Figure 1c.

Other contributors to the drop in performance is the global synchronizations and false sharing occurring while executing the kernel. These factors strongly depend on the number of colors and in general increase with colors required to resolve the dependency. A rough estimate for the barrier/synchronization overhead can be evaluated by comparing performance between kernel with (correct result) and without barriers (incorrect result) as seen in Figure 1d. Figure 1c shows the volume of false-sharing as measured by LIKWID for SpMTV kernels. For SpMV the amount of false sharing is negligible and is close to zero since it does not have indirect writes. It has to be noted that LIKWID measurements for counters like false sharing, might not be 100% accurate, but they provide a good qualitative estimate for comparison, as seen here counter indicates increase in false sharing with number of threads. More on counter validation can be found in [20].

It was seen that for most of the matrices arising in the project the average drop in performance by multi-coloring was almost a factor of two on a single socket of Ivy-Bridge. Although for most of the matrices the performance could be improved by algebraic block multi-coloring (ABMC), still the results we obtained were not optimal mainly for large matrices when compared to performance models which we



will see later in section 7. This led to the development of a method which works on a common data format like CRS in which most of the other kernels are written and at the same time preserves data locality, reduce synchronization overheads and false sharing.

**6. RACE method.** Keeping in mind the observations from previous section 5, one could observe that it would be best to maintain the non-zeros of matrix close to the diagonal. This has been observed previously in the regard of normal sparse matrix computations like SpMV and has led to the pre-processing of matrix by applying bandwidth reduction algorithms like “Reverse Cuthill McKee ” (RCM ). Now we aim to develop a method that does not distort this ideal permutations to a large extent but at the same time resolve distance- $k$  dependencies.

Our approach can be seen as a recursive level based method. Each step of the method basically consists of four steps namely:

1. Level construction
2. Permutation
3. Distance- $k$  coloring
4. Load balancing

The method is strongly coupled to the hardware underneath and exploits only the parallelism as required by the hardware. If at the end of all these four steps one does not achieve sufficient parallelism, all the steps are recursively applied to selected sub-graphs of the matrix until sufficient parallelism is attained. This recursive nature of our coloring method led to the naming of the method as “Recursive Algebraic Coloring Engine ” in short RACE .

To explain the method in an easier and illustrative way we choose a simple matrix namely the 2D 7pt. stencil. The sparsity pattern and the corresponding graph of the matrix is as shown in Figure 3.

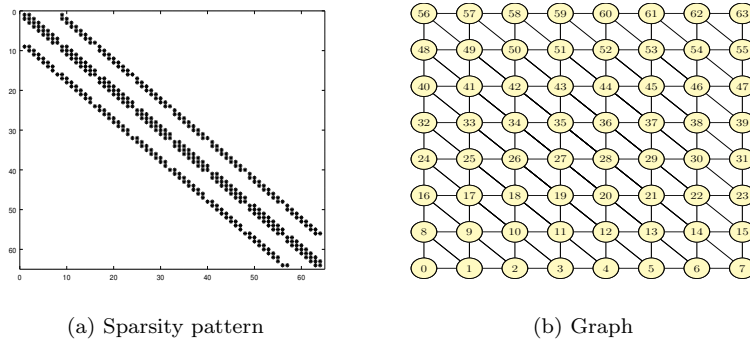


FIG. 3. 2d-7pt Stencil

**Definitions.** The following basic definitions from graph theory are used in the following sections:

- **Graph :**  $G = (V, E)$  represents a graph where  $V(G)$  belongs to set of vertices and  $E(G)$  represents the edges in the graph. Note that here we specifically denote  $G$  for fully connected undirected graphs.
- **Neighborhood :** Neighborhood of vertex  $u$  represented as  $N(u)$  is defined

as:

$$N(u) = \{v : uv \in E\}$$

- **Subgraph** : A subgraph  $H$  of graph  $G$  in this paper specifically refers to subgraph induced by  $V' \subseteq V(G)$  and is defined as

$$H = (V', \{uv : uv \in E(G) \text{ and } u, v \in V'\})$$

**6.1. Level Construction.** The first step of the RACE method is level construction. The step concerns with finding different *levels* in the graph. First *level* ( $L(0)$ ) is chosen to consist of a selected a root vertex. Level  $i$  ( $L(i)$ ) for  $i > 0$  is defined to contain vertices that are in neighborhood of vertices in previous *level* ( $L(i-1)$ ) and not in  $L(i-2)$  i.e.,

$$L(i) = \begin{cases} u : u \in N(L(i-1)) \cap N(L(i-2))' & \text{if } i \neq 0 \\ \text{root} & \text{otherwise} \end{cases}$$

These *levels* are same to the ones found in “Breadth First Search” (BFS) traversal and aims to find the minimum distance from the root node. Note that this is substantially different to the *levels* in methods like “level-scheduling” [21] where maximum distance is sought after. Algorithm A.1 shows an algorithm to find each nodes minimum distance from root,  $L(i)$  consists of all the nodes having a distance of  $i$  from the root. Total number of levels obtained with this graph traversal will be denoted as *total\_level*. Figure 4a shows *levels* on the 2d-7pt stencil (*total\_level* = 14), the main number on each vertex refers to the vertex number and the superscript shows the *level* number, i.e.,

$$n^i \implies n \in L(i)$$

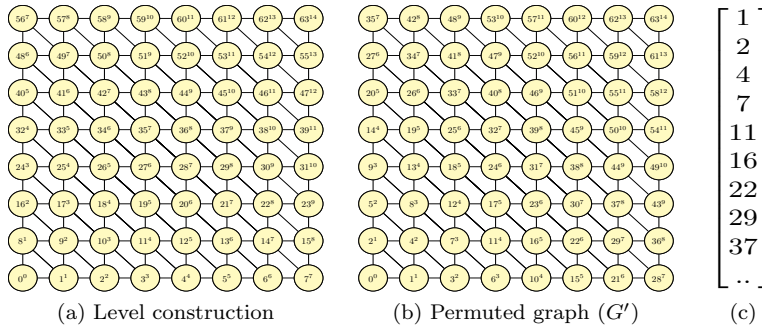


FIG. 4. (a) Levels in 2d-7pt stencil, (b) shows graph  $G'$  after permutation and (c) is the associated `level_ptr` to  $G'$ .

**6.2. Permutation.** Once the *levels* are known one has to permute the matrix in the order of its *levels*, such that vertices in  $L(i)$  appears before that of  $L(i-1)$ . Till this step the procedure is similar to that of BFS pre-processing for bandwidth reduction. One could also replace BFS with better bandwidth reduction algorithms like “(Reverse) Cuthill McKee”. Figure 4 shows the graph ( $G' = P(G)$ ) of 2d-7pt stencil matrix after this permutation ( $P$ ) is applied. Observe the difference in

node numbering between original lexicographic ordering in Figure 4a and Figure 4b. Now the most important step for resolving dependencies (coloring) is to store the information about *levels*. In order to do this we use a data structure called `level_ptr`. It stores the starting vector of each *levels*, which implies that *levels* on  $G'$  can be identified as:

$$L(i) = \{ u : u \in [\text{level\_ptr}[i] : (\text{level\_ptr}[i+1] - 1)] \text{ and } u \in V(G') \}$$

`level_ptr` for 2d-7pt stencil example is shown in Figure 4c, and one could easily read from `level_ptr` that vertices from `level_ptr(4) = 7` to `level_ptr(5) - 1 = 10` belongs to  $L(4)$ .

**6.3. Distance- $k$  coloring.** Two vertices are called distance- $k$  neighbours if the shortest path connecting them consists of at most  $k$  edges. This implies  $u$  is a distance- $k$  neighbour of  $v$  (denoted as  $u \xrightarrow{k} v$ ) if

$$(6.3) \quad u \xrightarrow{k} v \iff v \in \{ u \cup N(u) \cup N^2(u) \cup \dots \cup N^k(u) \}$$

Since we consider only undirected graph  $u \xrightarrow{k} v$  also implies  $v \xrightarrow{k} u$ . After having the permuted graph  $G'$  one can show that  $L(i)$  and  $L(i+k+j)$  where  $j \geq 1$  are distance- $k$  independent as shown in the following Corollary 6.1:

COROLLARY 6.1.  $L(i)$  and  $L(i \pm (k+j))$  are distance- $k$  independent  $\forall j \geq 1$ .

*Proof.* We prove by contradiction. Let there exist  $u, v \in V(G')$  such that  $u \in L(i)$  and  $v \in L(i \pm (k+j)) \forall j \geq 1$ . Assume  $u, v$  are distance- $k$  neighbours ( $u \xrightarrow{k} v$ ). From (6.1), (6.3) and the fact  $G'$  is undirected we get

$$\begin{aligned} u \xrightarrow{k} v &\iff v \in \{ L(i) \cup L(i \pm 1) \cup \dots \cup L(i \pm k) \} \\ &\implies v \notin L(i \pm (k+j)) \forall j \geq 1 \end{aligned}$$

which is a contradiction to the fact  $v \in L(i \pm (k+j)) \forall j \geq 1$ , this implies  $u$  and  $v$  are distance- $k$  independent.  $\square$

Corollary 6.1 implies that if we leave a gap of *at least* one *level* between any two *levels* ( $L(i), L(i+2)$  for example) all the vertices between them are distance-1 independent. Similarly if there is a gap of *at least* two *levels* between any two *levels* ( $L(i), L(i+3)$  for example) we get distance-2 independent levels.

Due to this weak definition in Corollary 6.1 there exists many possibility to make *levels* independent of each other and Figure 5 shows one such possibility each for distance-1 and distance-2 independent *levels*. One could group some of the nearby *levels* together to form a *level group*, and make this distance-1 or distance-2 independent of other *level groups*. The  $i$ -th *level group* would be denoted by  $T(i)$ . Difference between *level* and *level group* can be seen in Figure 5b, for Figure 5a *level group* and *level* coincides.

In principle one could compute on all independent *level groups* in parallel, but serial within a *level group*, i.e. for example in Figure 5b  $T(0), T(2), T(4), T(6)$  can be operated by four different threads in parallel and in the next sweep rest *level groups*. For the configurations seen in Figure 5 this would mean we have  $\frac{\text{total\_level}}{2}$  and  $\frac{\text{total\_level}}{4}$  parallelism for distance-1 and distance-2 kernels respectively.

But the problem with the configurations like the one seen in Figure 5 is that there is load imbalances between threads as the non-zeros per *level group* is not distributed

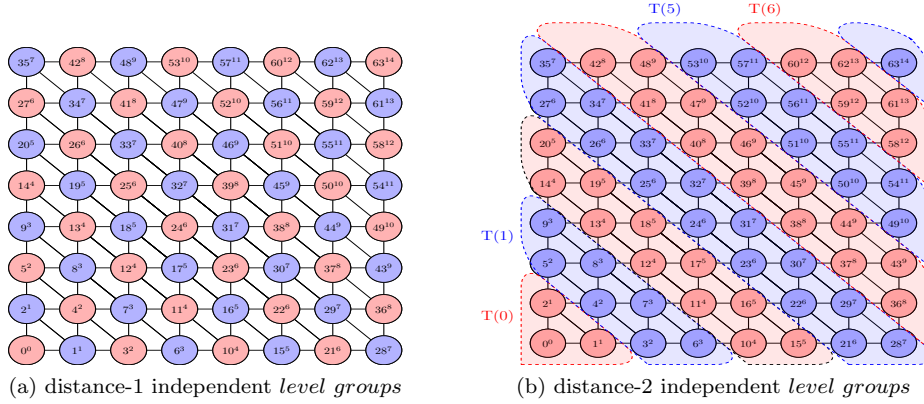


FIG. 5. distance-1 and distance-2 independent level groups .

evenly. As seen here in the case of 2d-7pt stencil the threads working on extreme ends of graph (e.g.,  $T(1), T(7)$ ) have small amount of work compared to the threads working on middle (e.g.,  $T(3), T(4)$ ).

**6.4. Load balancing.** Depending on the matrix each *level group* would contain different number of non-zeros, which leads to load imbalances as seen above in subsection 6.3. In order to avoid this problem we employ a load balancing scheme. At this step we plug in detail from hardware side like total parallelism. The idea is to exploit only the parallelism as required by the hardware while at the same time maintain distance- $k$  constraint seen in Corollary 6.1. To balance the load more nearby *levels* would be added to a *level group* which has less number of non-zeros and at *level group* where we have considerably big *levels* only sufficient amount of *levels* to maintain distance- $k$  constraint would be assigned. Assigning nearby levels instead of a random level further helps in preserving data locality.

An algorithm for load balancing can be found in Algorithm A.2. The aim of the algorithm is to reduce combined variance of number of non-zeros in each *level group* ( $T\_size(i)$  refers to non-zeros in *level group*  $i$ ). It does this by calculating mean and variance of  $T\_size$  in each parallel sweeps. For example in Figure 5b we need to calculate mean of  $T\_size$  of all *level groups* in red sweep and blue sweep separately. The combined variance is then found by summing up the variances in each parallel sweep. In order to reduce this combined variance we select the *level group* that has biggest absolute deviation from mean and try to add/remove levels to/from this *level group* from/to a *level group* that has biggest/least signed deviation. While removing *levels* from a *level group* one has to take care that the distance- $k$  coloring is not violated, for example in case of distance-2 and two sweep scheme like seen in Figure 5b we need to ensure at least two levels remain in a *level group*. To aid this shifting of *levels* to/from *level group* we use the pointers to *level group* denoted by  $T\_ptr$ . Doing this process in an iterative way finally we end up in a state with lowest combined variance at which no further moves are possible either due to violation of distance- $k$  dependency or due to increase in combined variance. Figure 6 shows step by step procedure involved in load balancing and Figure 7 shows *level groups* after load balancing applied on 2d-7pt stencil example of size  $16 \times 16$ .

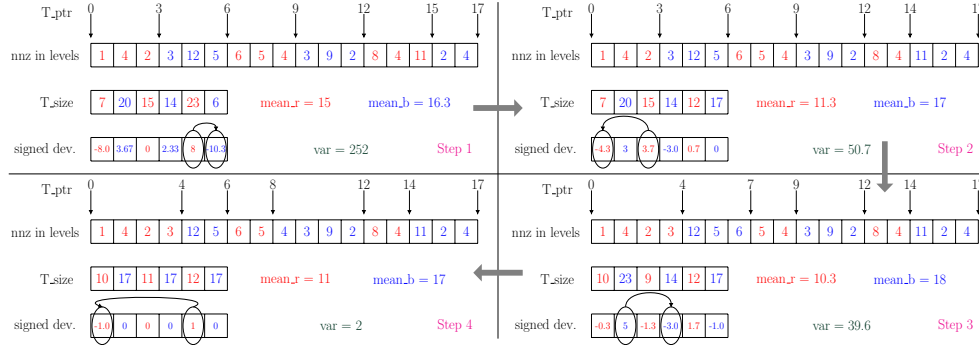
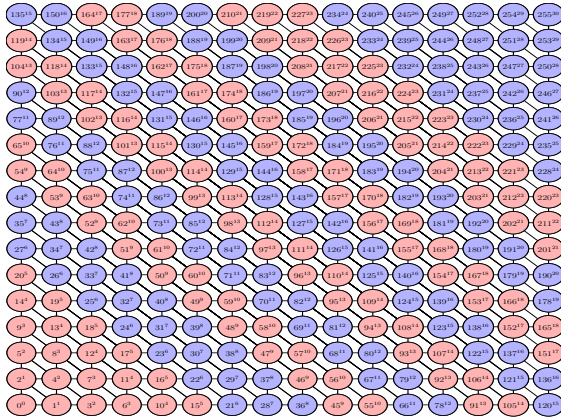


FIG. 6. Steps in load balancing (clockwise starting from top-left)

FIG. 7. After load balancing for five threads and distance-2 dependency on 2d-7pt stencil example, domain size  $16 \times 16$ . Note that level groups at extreme end have more levels due to less number of non-zeros in each level, while level groups in middle having bigger levels maintain two levels to preserve distance-2 constraint.

**6.5. Recursion.** As seen above in subsection 6.3 maximum amount of parallelism by the above approach depends on *total\_level*, also for most of the graphs as we approach the limit of parallelism there is not much room for load balancing, leading to imbalances. Depending on matrix and hardware underneath this might lead to inefficient utilization of resources. In order to avoid this problem we use the concept of recursion and exploit further parallelism if required by the hardware. Idea here is to intelligently select sub-graph(s) of the entire matrix and apply all the four steps recursively on this sub-graph. In the following we will show this concept in the context of distance-1 and later we will extent it to distance- $k$  dependencies. Further we will discuss on the method employed to select proper sub-graph and to have a globally balanced load.

**6.5.1. Distance-1.** *Level groups* which we constructed till now belongs to stage 1 of recursion and to make the explanations easier the stage number of recursion would be denoted as subscript i.e.,  $L_s(i)$  denotes *level i* of stage  $s$ . Contrary to methods like multi-coloring we didn't require each nodes in a color to be distance-1 independent of each other rather we had a weak constraint as prescribed by Corollary 6.1. Due to this there can exist more parallelism within a *level group*. For example in Figure 8 we see that within third *level group* ( $T_1(3)=L_1(3)$ ) vertices  $4 \not\rightarrow 5$  (4 distance-1 independent to 5),  $4 \not\rightarrow 6$ ,  $4 \not\rightarrow 7$  and  $5 \not\rightarrow 7$ , implying each of these pairs can be

388 computed in parallel without any distance-1 conflicts. This parallelism couldn't be  
 389 exploited in stage 1 since vertices in  $L_1(k)$  (here  $k=3$ ) were connected to preceding  
 390 level  $L_1(k-1)$  although some of them were not distance-1 dependent within  $L_1(k)$ .  
 391 In order to exploit this parallelism we use the concept of recursion.

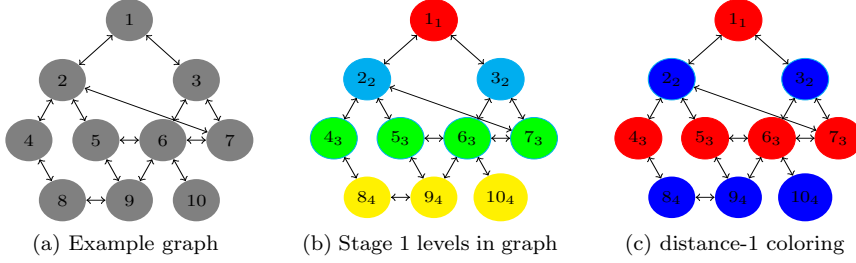


FIG. 8. Shows potential for more parallelism.  $T_1(2)$ ,  $T_1(3)$  and  $T_1(4)$  has more parallelism.

392 Recursion begins by selection of a sub-graph of the matrix. A typical choice is a  
 393 sub-graph induced by vertices in a *level group* of previous stage, more on the selection  
 394 of sub-graph will be seen later in subsection 6.5.4. For example let's choose sub-graph  
 395 induced by  $T_1(3)$  for recursion. The chosen sub-graph can be isolated from rest of  
 396 the graph since distance-1 coloring step in stage 1 has already made *level groups* in  
 397 a sweep independent of each other. Now we just need to repeat all the four step  
 398 explained previously (subsection 6.1 - subsection 6.4) to exploit parallelism within  
 399 this sub-graph .

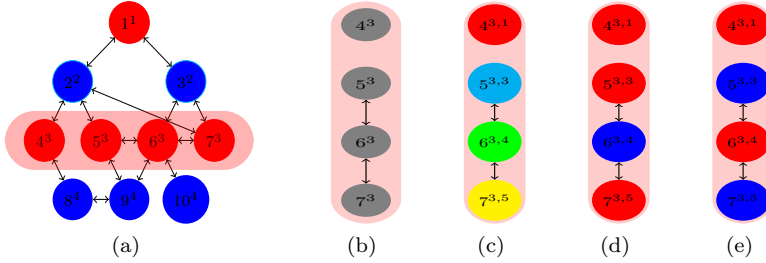


FIG. 9. Shows recursion being applied to  $T_1(3)$ . Figure 9b shows the selected sub-graph ,  
 Figure 9c shows level construction step on the sub-graph , Figures 9d and 9e shows two possibility  
 of distance-1 coloring of the sub-graph

400 Figure 9 shows an illustration of applying stage 2 of recursion on  $T_1(3)$  to find  
 401 more parallelism. To incorporate the information of levels after recursion we extent  
 402 the definition in (6.2) to the following:

$$(6.4) \quad n^{i,j,k,\dots} \implies n \in \{L_1(i) \cap L_2(j) \cap L_3(k) \cap \dots\}$$

404 Note that the sub-graphs might have multiple islands (group of vertices in a graph  
 405 that are not connected to rest of the graph). For example vertex 4 in Figure 9b  
 406 is an island in the considered sub-graph , similarly vertices 5,6,7 combine to form an  
 407 island. Since an island is totally disconnected from the rest of the graph it can be  
 408 executed in parallel to rest of the graph. To take advantage of this the starting node  
 409 in next island is assigned with an increment of two levels, as seen in Figure 9c. Due to



this there exists multiple valid distance-1 configuration (here Figures 9d and 9e) and the selection of the optimum one will be done in the final load balancing step of a particular stage as described in subsection 6.4.

With this recursive process we were able to find independent *level groups* ( $T_{s+1}$ ) within *level group* of previous stage ( $T_s$ ) and therefore the thread which works on  $T_s$  has to spawn threads to parallelize within  $T_{s+1}$ .

**6.5.2. Distance- $k$ .** For distance- $k$  the same procedure as distance-1 applies, except with a slight difference in selecting the sub-graph. In distance-1 we considered sub-graphs induced by *level groups*, but for distance- $k$  coloring this is not sufficient. As seen in Figure 10 for distance-2 coloring the selection of  $T_1(2)$  as sub-graph did not guarantee distance-2 independency between *level group*  $T_2$  within the sub-graph. This is due to the fact for  $k > 1$  dependency vertices  $a, b$  within a sub-graph might be connected to a common vertex ( $c$ ) outside the sub-graph leading to a distance- $k$  dependency between  $a$  and  $b$ . In Figure 10 we see  $4 \xrightarrow{1} 2$  &  $7 \xrightarrow{1} 2 \implies 4 \xrightarrow{2} 7$ , but since vertex 2 was not in the sub-graph considered we missed this dependency.

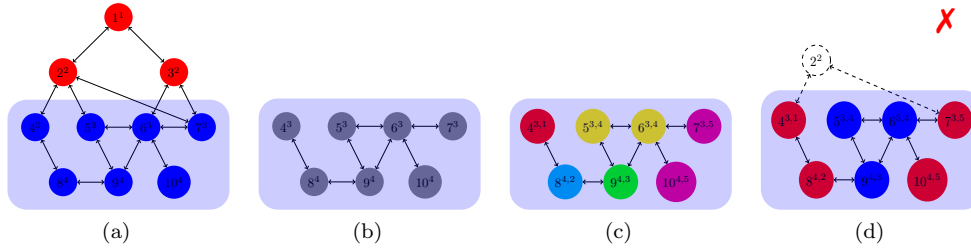


FIG. 10. Figures 10a and 10b shows level group induced sub-graph selected for recursion in case of distance-2. But applying the four steps to this selected sub-graph does not guarantee a distance-2 independency between level group of same sweep (color) as seen in Figure 10d

In order to resolve such dependency we have to consider an extra  $(k-1)^{th}$  interface level(s) of the selected sub-graph for the level construction step.  $k^{th}$  interface level of subgraph  $L_s(j)$ , denoted as  $I^k(L_s(j))$ , is defined as follows:

$$I^k(L_s(j)) = \{ u : u \xrightarrow{k} v \forall v \in L_s(j) \text{ and } u \notin L_s(j) \}$$

For distance-2 this would mean we have to include 1 interface level, the new selection is illustrated in Figure 11. With the new sub-graph selection for distance-2 coloring as seen in Figure 11a, the result after third step distance- $k$  coloring remains correct. In the example vertices 4 and 7 which had same color previously now gets a different color in (see Figure 11d).

Note that the interface levels have to be considered only in the first step namely level construction in the rest of the steps we just need to consider target sub-graphs induced by *level groups*.

**6.5.3. level\_tree.** By recursion we are able to exploit more parallelism. However this introduces more complexity and one has to also respect the dependencies between stages in addition to one within stages. The best idea is to have a data structure similar to the recursion, therefore we extent the `level_ptr` data structure to a hierarchical tree data structure to store the informations. This data structure is called a `level_tree`. The root of `level_tree` contains information of entire domain,

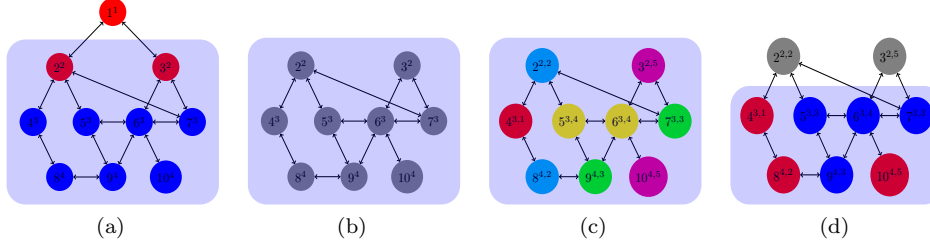


FIG. 11. Correct procedure of selecting sub-graph for distance-2 coloring. The level group  $T(2)$  and its 1<sup>st</sup> interface level is chosen as shown in Figures 11a and 11b for level construction stage seen in Figure 11c. For rest of the steps only required sub-graph to be parallelised is considered as shown in Figure 11d for distance-k coloring.

first leaves of this root stores information about *level groups* in stage 1 ( $T_1$ ), next  
leaves about *level groups* in stage 2 ( $T_2$ ) and so on.

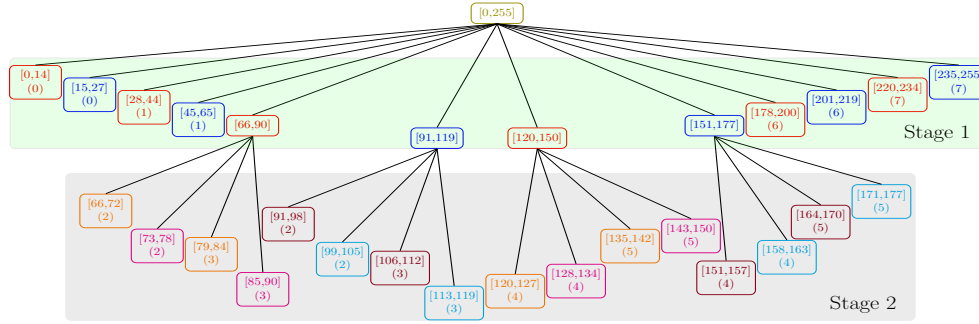
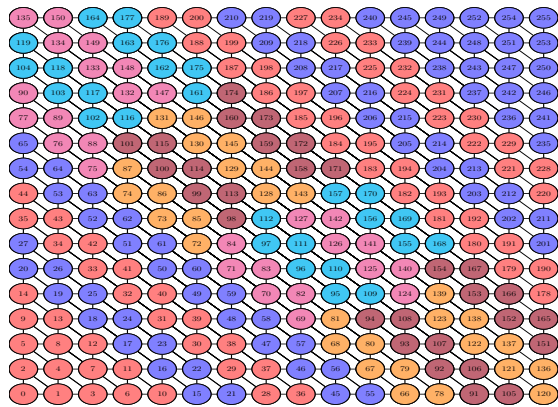


FIG. 12. **leveltree** corresponding to 2d-7pt stencil example for domain size  $16 \times 16$ , and 8 threads. The range specified in each leaves represent the vertices belonging to each level group, and the number in bracket represents the thread assigned to the level group in fill type pinning.



for parallel all red  
for parallel all orange  
for parallel all pink  
for parallel all blue  
for parallel all brown  
for parallel all cyan

FIG. 13. Graph corresponding to **leveltree** in Figure 12. The execution order of different level group is specified in the short code snippet on right. Note nested parallelism being used.

Figure 12 shows a **leveltree** corresponding for 2d-7pt stencil example. Threads  
are assigned to each *level group*, order of which depends on the pinning strategy  
used. For example in fill type pinning strategy one would pin thread 0 to  $T_1(0)$   
and  $T_1(1)$ , thread 1 to  $T_1(2)$  and  $T_1(3)$ , thread 2 to  $T_2(0) \subset T_1(4)$ ,  $T_2(1) \subset T_1(4)$ ,



$T_2(0) \subset T_1(5)$  and  $T_2(1) \subset T_1(5)$ , and so on. In order to replicate this tree like parallelisation strategy we use nested parallelism, where threads in stage  $k + 1$  is spawned by threads in stage  $k$ . The graph corresponding to 2d-7pt stencil example is shown in Figure 13, and the execution order is specified in the figure. At the end of each `for parallel all color` there is synchronization between threads assigned to *level group* of corresponding color. Since each of the leaf need to synchronize only with it's siblings (leaves of same parent) we use simple point to point synchronization scheme.

**6.5.4. Sub-graph selection and global load balancing.** Parallelism required for hardware underneath can be obtained either by expanding the `level_tree` horizontally i.e., increasing *level groups* within a stage or by expanding `level_tree` vertically with the help of recursion. But as we have seen before in subsection 6.3 the horizontal parallelism is limited and after a certain extent this would lead to load balancing. Similarly excessive usage of recursion is also not a good idea since data locality worsens due to local permutations within sub-graph. Therefore it is vital to find a proper balance and choose proper configuration. Furthermore just doing load balancing within a single stage is not the best, for example if we had equally balanced within stage 1 in Figure 12, we would receive no benefit from recursion. Therefore a global load balancing becomes inevitable.

In order to select proper sub-graph and do global load balancing we employ a simple algorithm to find proper weights for each *level group* ( $T_s(i)$ ) in a particular stage, then depending on this weights, denoted as  $w(T_s(i))$ , we do load balancing with weights in the particular stage (as seen in Algorithm A.2, except weightage is given to *level groups*). Finally if  $w(T_s(i)) > 1$  we use recursion to achieve  $w(T_s(i))$  parallel work in the next stage of  $T_s(i)$ . The basic structure of the algorithm employed to find weights is as follows:

1. Find weights,  $w(L_s(i))$  for each level in the current stage ( $s$ ) by

$$w(L_s(i)) = (\text{level\_ptr}_s[i + 1] - \text{level\_ptr}_s[i]) * \frac{n\_threads}{n\_vertices}$$

$n\_threads$  : total parallelism required by hardware

$n\_vertices$  : number of vertices in graph

2. Starting from  $w(L_s(0))$  sum up weights till they form a number ( $a$ ) close to whole number ( $b$ ). The closeness can be controlled by an efficiency parameter for stage  $s$ ,  $\epsilon_s$  is defined as:

$$\epsilon_s = 1 - \text{abs}(a - b);$$

The obtained number  $b$  is chosen as weight for *level groups* operated by first thread in the current stage i.e.,  $w(T_s(0)) = w(T_s(1)) = b$ . The weight for next *level groups* are found by resetting the sum counter to zero and repeating the procedure with *levels* just after the current *level group*.

**7. Experiments and Results.** The method stated above was implemented and consolidated into a library called RACE. The library provides easy interface for parallelizing kernels, user typically just needs to supply the serial code (with dependency) and hardware settings. Library will then parallelize, pin and run the code in parallel. The library is publicly available in the git repository.

In the following we present the performance and convergence results obtained

494 using the library, and compare it against state of art methods. Test setup, hardware  
 495 and matrices as described in section 4 is used for the following benchmarks.

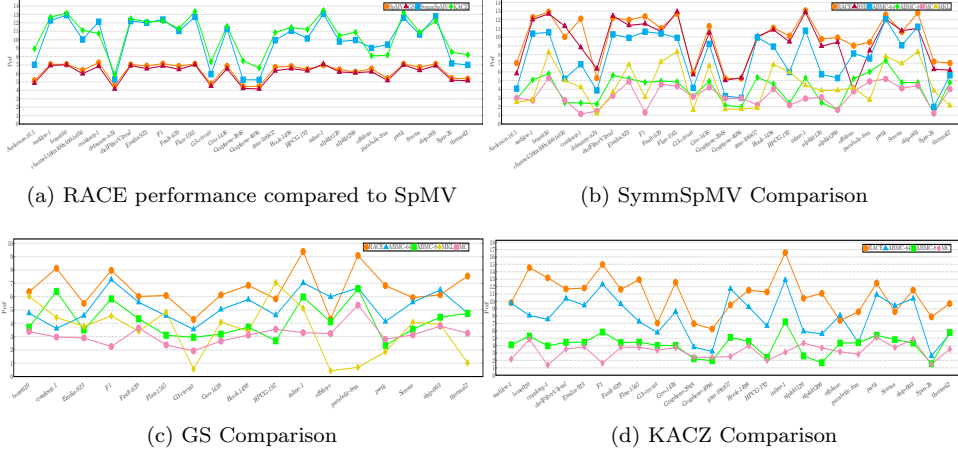


FIG. 14. Performance results on Ivy-Bridge

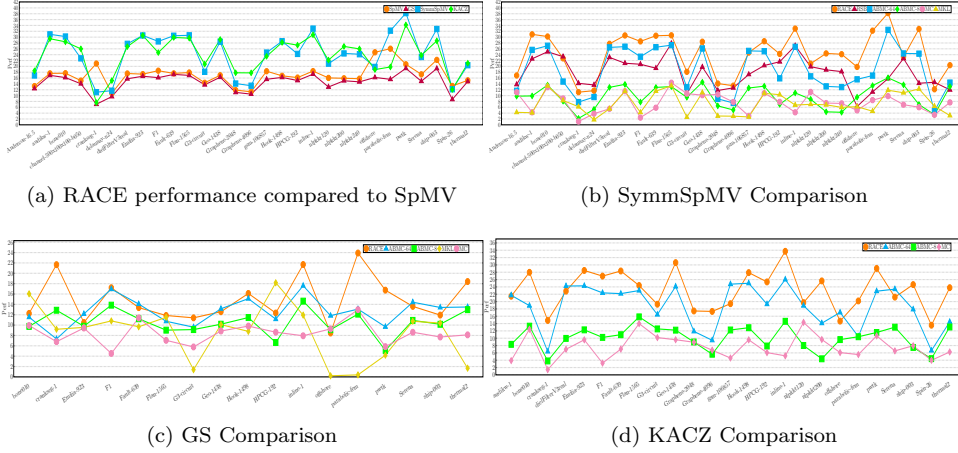


FIG. 15. Performance results on Sky-Lake

### 7.1. Main points to discuss.

- Relate roofline model and the performance graphs of RACE compared to SpMV.
- Point out on Ivy-Bridge we reach close to ideal performance in every case, and on Sky-Lake except for corner cases like crankseg and offshore we reach close to ideal performance. The drop in corner cases like crankseg and offshore on Sky-Lake is due to lack of parallelism attained by RACE and associated load imbalances. This effect shows up on Sky-Lake rather than Ivy-Bridge since Sky-Lake has 20 threads compared to 10 on Ivy-Bridge .
- Point out that for cases like Graphene, Spin, parabolic\_fem we don't see 2 fold increase in GFlop/s for KACZ, and SymmSpMV. This is due to the

fact here  $N_{n_{zr}}$  is very small like 4, 14 and 7 which causes two problems. For KACZ kernel there is one division per row and this causes a performance drop as evident in Spin matrices, also this effect can be observed for GS kernel. For SymmSpMV kernel the  $N_{n_{zr}}$  decreases almost by half since we operate only on upper triangular part and with short loop over  $N_{n_{zr}}$  no effective vectorization and modulo unrolling can be done.

- Matrices like crankseg-1, and offshore are also really small making some part of data fit in cache, this is the reason why they achieve performance above RLM.
- Discuss why we chose the methods for comparison. MC and ABMC are common in literature for distance-1 coloring, MKL methods are standard library used in many productive codes, also it uses level-scheduling (not explicitly stated but we believe) for kernels like GS and enables us to compare with methods that do not disturb convergence. RSB enables to compare with methods using different data format and it has been shown this method has an upper hand in this category.
- Comparison with SymmSpMV shows the behavior of different methods for distance-2 coloring. Here we see in almost all of the case RACE and RSB has an upper hand on Ivy-Bridge . Although in some cases like offshore RACE clearly has an advantage. ABMC methods follow these methods. MKL and MC does not deliver good performance. For Sky-Lake architecture RSB falls behind ABMC, we think this is because of the requirement of ABMC to lock rows and cols of the submatrix on which a thread is working, becoming a bottleneck at high thread counts.
- Maybe tell RSB and 16-bit integer.
- Discuss with methods like ABMC and MC the performance especially drops for large matrices like Graphene, Spin, nlpkkt due to worsening of data locality ( $\alpha$ ). Show sparsity pattern and LIKWID measurements.
- Tell GS and KACZ performance includes also takes iterations into consideration (as shown in paper). Tell we do only a distance-1 coloring for GS and distance-2 for KACZ. We use only matrices where GS can be applied and similarly for KACZ. Also we just compare against readily available solutions. Therefore RSB is left out for GS and RSB and MKL left out for KACZ.
- For GS RACE has an upper hand on Ivy-Bridge and on Sky-Lake RACE and ABMC have almost similar performance on Sky-Lake , although for some cases RACE has huge advantage. Reason for this advantage is due to slight decrease in iterations for RACE (put fig) and slight improvement in performance compared to ABMC for distance-1 case. For offshore case RACE performs worse than ABMC, this is because here with RACE one requires more iterations. Also note that all the large matrices which we had are unsuitable for GS sweep as they do not converge, but for large matrices the performance drops again for ABMC method due to degrading of  $\alpha$  factor. (Maybe just put perf. pictures).
- Main advantage of RACE method comes with kernels having distance-2 dependencies like SymmSpMV and KACZ since here methods like ABMC require more colors and their locality degrades further since here within a color rows have to be structurally orthogonal (rows shouldn't have common column entries). Performance on KACZ shows this advantage. Here we again see for moderately large matrices the advantage is higher. Iteration behavior between methods remains similar to GS.

557     **8. Application Runs.**

558     **9. Conclusion.**

559     **10. Future Work.**

560     **Acknowledgments.** We would like to acknowledge the assistance of volunteers  
561 in putting together this example manuscript and supplement.

## REFERENCES

- [1] *Equipping Sparse Solvers for Exascale - ESSEX*. <https://blogs.fau.de/essex/activities>.
- [2] *SuiteSparse Matrix Collection*. <https://sparse.tamu.edu/>.
- [3] A. BULUÇ, J. T. FINEMAN, M. FRIGO, J. R. GILBERT, AND C. E. LEISERSON, *Parallel sparse matrix-vector and matrix-transpose-vector multiplication using compressed sparse blocks*, in Proceedings of the Twenty-first Annual Symposium on Parallelism in Algorithms and Architectures, SPAA '09, New York, NY, USA, 2009, ACM, pp. 233–244, <https://doi.org/10.1145/1583991.1584053>.
- [4] J. DONGARRA AND M. HEROUX, *Toward a new metric for ranking high performance computing systems*, Tech. Report SAND2013-4744, Sandia National Laboratories, 2013.
- [5] T. ELFVING, *Block-iterative methods for consistent and inconsistent linear equations*, Numerische Mathematik, 35 (1980), pp. 1–12, <https://doi.org/10.1007/BF01396365>, <https://doi.org/10.1007/BF01396365>.
- [6] D. J. EVANS, *Parallel S.O.R. Iterative Methods*, Parallel Comput., 1 (1984), pp. 3–18, [https://doi.org/10.1016/S0167-8191\(84\)90380-6](https://doi.org/10.1016/S0167-8191(84)90380-6).
- [7] M. GALGON, L. KRÄMER, J. THIES, A. BASERMANN, AND B. LANG, *On the parallel iterative solution of linear systems arising in the fast algorithm for computing inner eigenvalues*, Parallel Comput., 49 (2015), pp. 153–163, <https://doi.org/10.1016/j.parco.2015.06.005>.
- [8] T. GKOUNTOUVAS, V. KARAKASIS, K. KOURTIS, G. GOUMAS, AND N. KOZIRIS, *Improving the performance of the symmetric sparse matrix-vector multiplication in multicore*, in 2013 IEEE 27th International Symposium on Parallel and Distributed Processing, May 2013, pp. 273–283, <https://doi.org/10.1109/IPDPS.2013.43>.
- [9] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, The Johns Hopkins University Press, Baltimore, 4th ed., 2013.
- [10] D. GORDON AND R. GORDON, *Component-averaged row projections: A robust, block-parallel scheme for sparse linear systems*, SIAM Journal on Scientific Computing, 27 (2005), pp. 1092–1117, <https://doi.org/10.1137/040609458>.
- [11] T. IWASHITA, H. NAKASHIMA, AND Y. TAKAHASHI, *Algebraic block multi-color ordering method for parallel multi-threaded sparse triangular solver in iccg method*, in Proceedings of the 2012 IEEE 26th International Parallel and Distributed Processing Symposium, IPDPS '12, Washington, DC, USA, 2012, IEEE Computer Society, pp. 474–483, <https://doi.org/10.1109/IPDPS.2012.51>, <http://dx.doi.org/10.1109/IPDPS.2012.51>.
- [12] M. T. JONES AND P. E. PLASSMANN, *Scalable iterative solution of sparse linear systems*, Parallel Comput., 20 (1994), pp. 753–773, [https://doi.org/10.1016/0167-8191\(94\)90004-3](https://doi.org/10.1016/0167-8191(94)90004-3), [http://dx.doi.org/10.1016/0167-8191\(94\)90004-3](http://dx.doi.org/10.1016/0167-8191(94)90004-3).
- [13] C. KAMATH AND A. SAMEH, *A projection method for solving nonsymmetric linear systems on multiprocessors*, 9 (1989), pp. 291–312, [https://doi.org/10.1016/0167-8191\(89\)90114-2](https://doi.org/10.1016/0167-8191(89)90114-2).
- [14] M. KREUTZER, G. HAGER, G. WELLEIN, H. FEHSKE, AND A. R. BISHOP, *A Unified Sparse Matrix Data Format for Efficient General Sparse Matrix-Vector Multiplication on Modern Processors with Wide SIMD Units*, SIAM Journal on Scientific Computing, 36 (2014), pp. C401–C423, <https://doi.org/10.1137/130930352>.
- [15] H. LU, M. HALAPPANAVAR, D. CHAVARRA-MIRANDA, A. H. GEBREMEDHIN, A. PANYALA, AND A. KALYANARAMAN, *Algorithms for balanced graph colorings with applications in parallel computing*, IEEE Transactions on Parallel and Distributed Systems, 28 (2017), pp. 1240–1256, <https://doi.org/10.1109/TPDS.2016.2620142>.
- [16] M. MARTONE, *Efficient multithreaded untransposed, transposed or symmetric sparse matrix-vector multiplication with the recursive sparse blocks format*, Parallel Comput., 40 (2014), pp. 251–270, <https://doi.org/10.1016/j.parco.2014.03.008>, <http://dx.doi.org/10.1016/j.parco.2014.03.008>.
- [17] K. NAKAJIMA AND H. OKUDA, *Parallel iterative solvers for unstructured grids using an openmp/mpi hybrid programming model for the geofem platform on smp cluster architectures*, in High Performance Computing, H. P. Zima, K. Joe, M. Sato, Y. Seo, and M. Shimasaki, eds., Berlin, Heidelberg, 2002, Springer Berlin Heidelberg, pp. 437–448.
- [18] J. PARK, M. SMELYANSKIY, N. SUNDARAM, AND P. DUBEY, *Sparsifying synchronization for high-performance shared-memory sparse triangular solver*, in Proceedings of the 29th International Conference on Supercomputing - Volume 8488, ISC 2014, New York, NY, USA, 2014, Springer-Verlag New York, Inc., pp. 124–140, [https://doi.org/10.1007/978-3-319-07518-1\\_8](https://doi.org/10.1007/978-3-319-07518-1_8).
- [19] E. POLIZZI, *A density matrix-based algorithm for solving eigenvalue problems*, CoRR, abs/0901.2665 (2009), <https://arxiv.org/abs/0901.2665>.
- [20] T. RÖHL, J. EITZINGER, G. HAGER, AND G. WELLEIN, *Validation of hardware events for suc-*

- 623        *cessful performance pattern identification in high performance computing*, in Tools for High  
624        Performance Computing 2015, Cham, 2016, Springer International Publishing, pp. 17–28.
- 625        [21] Y. SAAD, *Iterative Methods for Sparse Linear Systems*, Society for Industrial and Ap-  
626        plied Mathematics, second ed., 2003, <https://doi.org/10.1137/1.9780898718003>, [https://](https://epubs.siam.org/doi/abs/10.1137/1.9780898718003)  
627        [epubs.siam.org/doi/abs/10.1137/1.9780898718003](https://epubs.siam.org/doi/abs/10.1137/1.9780898718003), <https://arxiv.org/abs/https://epubs.siam.org/doi/pdf/10.1137/1.9780898718003>.  
628        [siam.org/doi/pdf/10.1137/1.9780898718003](https://arxiv.org/abs/https://epubs.siam.org/doi/pdf/10.1137/1.9780898718003).
- 629        [22] J. TREIBIG, G. HAGER, AND G. WELLEIN, *Likwid: A lightweight performance-oriented tool*  
630        *suite for x86 multicore environments*, (2010).
- 631        [23] S. WILLIAMS, A. WATERMAN, AND D. PATTERSON, *Roofline: An Insightful Visual Performance*  
632        *Model for Multicore Architectures*, Commun. ACM, 52 (2009), pp. 65–76, [https://doi.org/](https://doi.org/10.1145/1498765.1498785)  
633        [10.1145/1498765.1498785](https://doi.org/10.1145/1498765.1498785).

## Appendix A. Algorithms.

---

### Algorithm A.1 Construction of levels

---

```

1: Choose starting node(s) =  $\{n\}$ 
2:  $marked\_all = \text{false}$ 
3:  $N = \text{nrows}(\text{graph})$ 
4:  $distFromRoot[1..N] = -1$ 
5:  $curr\_children.push\_back(n)$ ;
6:  $currLvl = 0$ 
7: while  $\neg marked\_all$  do
8:    $marked\_all = \text{true}$ 
9:    $nxt\_children = \{\}$ 
10:  for  $i = 1 : \text{size}(curr\_children)$  do
11:    if  $distFromRoot[curr\_children[i]] == -1$  then
12:       $distFromRoot[curr\_children[i]] = currLvl$ 
13:      for  $j$  in  $\text{graph}[curr\_children[i]].children$  do
14:        if  $distFromRoot[j] == -1$  then
15:           $nxt\_children.push\_back(j)$ 
16:        end if
17:      end for
18:    end if
19:  end for
20:   $curr\_children = nxt\_children$ 
21:   $currLvl = currLvl + 1$ 
22: end while

```

---

634

**Algorithm A.2** Load Balancing for two sweep, distance-2

---

```

1: num_sweep = 2                                % two sweep method
2: minGap = 2                                    %distance-2
3: len = num_sweep * nthread                    % constructing nthread parallel work
4: while !(exit) do
5:   T_size = update(T_ptr)                    %T_size contains non-zeros in each level group
6:   mean_r = sum(T_size[0 : num_sweep : len]) / nthreads
7:   mean_b = sum(T_size[1 : num_sweep : len]) / nthreads
8:   diff[0 : num_sweep : len] = T_size[0 : num_sweep : len]. - mean_r
9:   diff[1 : num_sweep : len] = T_size[1 : num_sweep : len]. - mean_b
10:  var = dot_product(diff, diff)
11:  absRankIdx = sortIdx(abs(diff)) % sortIdx returns permutation after
12:                                     % sorting from bigger to larger
13:  rankIdx = sortIdx(diff)
14:  currRank = 0, newVar = var
15:  old_T_ptr = T_ptr
16:  while newVar ≥ var do
17:    T_ptr = old_T_ptr
18:    fail=true
19:    if diff[absRankIdx[currRank]] < 0 then
20:      for el in rankIdx[(len - 1) : -1 : 0] do
21:        if (T_ptr[el + 1] - T_ptr[el]) > min_gap then
22:          acquireIdx = el
23:          fail=false
24:          break
25:        end if
26:      end for
27:      shift(T_ptr, acquireIdx, currRank) % shifts T_ptr by 1 from acquireIdx
28:                                     % to currRank if currIdx < acquireIdx else shift by -1
29:    else if (T_ptr[currRank + 1] - T_ptr[currRank]) > min_gap then
30:      giveIdx = rankIdx[0]
31:      fail=false
32:      shift(T_ptr, currRank, giveIdx)
33:    end if
34:    if !fail then
35:      newVar = calculate_variance(T_ptr) % as seen in Line 5 to Line 10
36:    end if
37:    if (currRank == (len - 1)) && (newVar ≥ var) then
38:      T_ptr = old_T_ptr
39:      exit = true
40:      break
41:    end if
42:    currRank += 1
43:  end while
44: end while

```

---