

# UC Berkeley

## UC Berkeley Previously Published Works

### Title

Exploiting Multiple Levels of Parallelism in Sparse Matrix-Matrix Multiplication

### Permalink

<https://escholarship.org/uc/item/9x19n6cd>

### Journal

SIAM Journal on Scientific Computing, 38(6)

### ISSN

1064-8275

### Authors

Azad, Ariful  
Ballard, Grey  
Buluç, Aydin  
[et al.](#)

### Publication Date

2016

### DOI

10.1137/15m104253x

Peer reviewed

# EXPLOITING MULTIPLE LEVELS OF PARALLELISM IN SPARSE MATRIX-MATRIX MULTIPLICATION

ARIFUL AZAD<sup>\*</sup>, GREY BALLARD<sup>†</sup>, AYDIN BULUÇ<sup>‡</sup>, JAMES DEMMEL<sup>§</sup>, LAURA GRIGORI<sup>¶</sup>, ODED SCHWARTZ<sup>||</sup>, SIVAN TOLEDO<sup>\*\*</sup>; AND SAMUEL WILLIAMS<sup>††</sup>

**Abstract.** Sparse matrix-matrix multiplication (or SpGEMM) is a key primitive for many high-performance graph algorithms as well as for some linear solvers, such as algebraic multigrid. The scaling of existing parallel implementations of SpGEMM is heavily bound by communication. Even though 3D (or 2.5D) algorithms have been proposed and theoretically analyzed in the flat MPI model on Erdős-Rényi matrices, those algorithms had not been implemented in practice and their complexities had not been analyzed for the general case. In this work, we present the first implementation of the 3D SpGEMM formulation that exploits multiple (intra-node and inter-node) levels of parallelism, achieving significant speedups over the state-of-the-art publicly available codes at all levels of concurrencies. We extensively evaluate our implementation and identify bottlenecks that should be subject to further research.

**Key words.** Parallel computing, numerical linear algebra, sparse matrix-matrix multiplication, 2.5D algorithms, 3D algorithms, multithreading, SpGEMM, 2D decomposition, graph algorithms.

**AMS subject classifications.** 05C50, 05C85, 65F50, 68W10

**1. Introduction.** Multiplication of two sparse matrices (SpGEMM) is a key operation for high-performance graph computations in the language of linear algebra [31, 40]. Examples include graph contraction [25], betweenness centrality [13], Markov clustering [47], peer pressure clustering [43], triangle counting [4], and cycle detection [49]. SpGEMM is also used in scientific computing. For instance, it is often a performance bottleneck for Algebraic Multigrid (AMG), where it is used in the set-up phase for restricting and interpolating matrices [7]. Schur complement methods in hybrid linear solvers [48] also require fast SpGEMM. In electronic structure calculations, linear-scaling methods exploit Kohn’s “nearsightedness” principle of electrons in many-atom systems [33]. SpGEMM and its approximate versions are often the workhorse of these computations [8, 10].

We describe new parallel implementations of the SpGEMM kernel, by exploiting multiple levels of parallelism. We provide the first complete implementation and large-scale results of a “3D algorithm” that asymptotically reduces communication costs compared to the state-of-the-art 2D algorithms. The name “3D” derives from the parallelization across all 3 dimensions of the iteration space. While 2D algorithms like Sparse SUMMA [14] are based on a 2D decomposition of the output matrix with computation following an “owner computes” rule, a 3D algorithm also parallelizes the computation of individual output matrix entries. Our 3D formulation relies on splitting (as opposed to replicating) input submatrices across processor layers.

While previous work [5] analyzed the communication costs of a large family of parallel SpGEMM algorithms and provided lower-bounds on random matrices, it did

<sup>\*</sup>CRD, Lawrence Berkeley National Laboratory, CA (azad@lbl.gov).

<sup>†</sup>Sandia National Laboratories, Livermore, California. Current address: Computer Science Department, Wake Forest University, Winston Salem, North Carolina (ballard@wfu.edu).

<sup>‡</sup>CRD, Lawrence Berkeley National Laboratory, CA (abuluc@lbl.gov).

<sup>§</sup>EECS, University of California, Berkeley, CA (demmel@eecs.berkeley.edu).

<sup>¶</sup>INRIA Paris-Rocquencourt, Alpines, France (laura.grigori@inria.fr).

<sup>||</sup>The Hebrew University, Israel (odedsc@cs.huji.ac.il).

<sup>\*\*</sup>Blavatnik School of Computer Science, Tel Aviv University, Israel (stoledo@tau.ac.il).

<sup>††</sup>CRD, Lawrence Berkeley National Laboratory, CA (swwilliams@lbl.gov).

not present any experimental results. In particular, the following questions were left unanswered:

- What is the effect of different communication patterns on relative scalability of these algorithms? The analysis was performed in terms of “the number of words moved per processor”, which did not take into account important factors such as network contention, use of collectives, the relative sizes of the communicators, etc.
- What is the effect of in-node multithreading? By intuition, one can expect a positive effect due to reduced network contention and automatic data aggregation as a result of in-node multithreading, but those have not been evaluated before.
- What is the role of local data structures and local algorithms? In particular, what is the right data structure to store local sparse matrices in order to multiply them fast using a single thread and multiple threads? How do we merge local triples efficiently during the reduction phases?
- How do the algorithms perform on real-world matrices, such as those with skewed degree distributions?

This paper addresses these questions by presenting the first implementation of the 3D SpGEMM formulation that exploits both the additional third processor grid dimension and the in-node multithreading aspect. In particular, we show that the third processor grid dimension navigates a tradeoff between communication of the *input* matrices and communication of the *output* matrix. We also show that in-node multithreading, with efficient shared-memory parallel kernels, can significantly enhance scalability. In terms of local data structures and algorithms, we use a priority queue to merge sparse vectors for in-node multithreading. This eliminates thread scaling bottlenecks which were due to asymptotically increased working set size as well as the need to modify the data structures for cache efficiency. To answer the last question, we benchmark our algorithms on real-world matrices coming from a variety of applications. Our extensive evaluation via large-scale experiments exposes bottlenecks and provides new avenues for research.

Section 3 summarizes earlier results on various parallel SpGEMM formulations. Section 4 presents the distributed-memory algorithms implemented for this work, as well as the local data structures and operations in detail. In particular, our new 3D algorithm, Split-3D-SpGEMM, is presented in Section 4.4. Section 5 gives an extensive performance evaluation of these implementations using large scale parallel experiments, including a performance comparison with similar primitives offered by other publicly available libraries such as Trilinos and Intel Math Kernel Library (MKL). Various implementation decisions and their effects on performance are also detailed.

**2. Notation.** Let  $\mathbf{A} \in \mathbb{S}^{m \times k}$  be a sparse rectangular matrix of elements from a semiring  $\mathbb{S}$ . We use  $nmz(\mathbf{A})$  to denote the number of nonzero elements in  $\mathbf{A}$ . When the matrix is clear from context, we drop the parenthesis and simply use  $nmz$ . For sparse matrix indexing, we use the convenient MATLAB colon notation, where  $\mathbf{A}(:, i)$  denotes the  $i$ th column,  $\mathbf{A}(i, :)$  denotes the  $i$ th row, and  $\mathbf{A}(i, j)$  denotes the element at the  $(i, j)$ th position of matrix  $\mathbf{A}$ . Array and vector indices are 1-based throughout this paper. The length of an array  $\mathbf{l}$ , denoted by  $len(\mathbf{l})$ , is equal to its number of elements. For one-dimensional arrays,  $\mathbf{l}(i)$  denotes the  $i$ th component of the array.  $\mathbf{l}(j : k)$  defines the range  $\mathbf{l}(j), \mathbf{l}(j + 1), \dots, \mathbf{l}(k)$  and is also applicable to matrices.

We use  $flops(\mathbf{A}, \mathbf{B})$ , pronounced “flops”, to denote the number of nonzero arithmetic operations required when computing the product of matrices  $\mathbf{A}$  and  $\mathbf{B}$ . When the operation and the operands are clear from context, we simply use  $flops$ . We ac-

knowledge that semiring operations do not have to be on floating-point numbers (e.g. they can be on integers or Booleans) but we nevertheless use flops as opposed to ops to be consistent with existing literature.

In our analysis of parallel running time, the latency of sending a message over the communication interconnect is  $\alpha$ , and the inverse bandwidth is  $\beta$ , both expressed as multiples of the time for a floating-point operation (also accounting for the cost of cache misses and memory indirections associated with that floating point operation). Notation  $f(x) = \Theta(g(x))$  means that  $f$  is bounded asymptotically by  $g$  both above and below. We index a 3D process grid with  $P(i, j, k)$ . Each 2D slice of this grid  $P(:, :, k)$  with the third dimension fixed is called a process “layer” and each 1D slice of this grid  $P(i, j, :)$  with the first two dimensions fixed is called a process “fiber”.

**3. Background and Related Work.** The classical serial SpGEMM algorithm for general sparse matrices was first described by Gustavson [27], and was subsequently used in Matlab [24] and CSparse [19]. For computing the product  $\mathbf{C} = \mathbf{AB}$ , where  $\mathbf{A} \in \mathbb{S}^{m \times l}$ ,  $\mathbf{B} \in \mathbb{S}^{l \times n}$  and  $\mathbf{C} \in \mathbb{S}^{m \times n}$ , Gustavson’s algorithm runs in  $O(\text{flops} + nnz + m + n)$  time, which is optimal when flops is larger than  $nnz$ ,  $m$ , and  $n$ . It uses the popular compressed sparse column (CSC) format for representing its sparse matrices. Algorithm 1 gives the pseudocode for this column-wise serial algorithm for SpGEMM.

---

**Algorithm 1** Column-wise formulation of serial matrix multiplication

---

```

1: procedure COLUMNWISE-SPGEMM( $\mathbf{A}, \mathbf{B}, \mathbf{C}$ )
2:   for  $k \leftarrow 1$  to  $n$  do
3:     for  $j$  where  $\mathbf{B}(j, k) \neq 0$  do
4:        $\mathbf{C}(:, k) \leftarrow \mathbf{C}(:, k) + \mathbf{A}(:, j) \cdot \mathbf{B}(j, k)$ 

```

---

McCourt et al. [41] target  $\mathbf{AB}^T$  and  $\mathbf{RAR}^T$  operations in the specific context of Algebraic Multigrid (AMG). A coloring of the output matrix  $\mathbf{C}$  finds *structurally orthogonal* columns that can be computed simultaneously. Two columns are structurally orthogonal if the inner product of their structures (to avoid numerical cancellation) is zero. They use matrix colorings to restructure  $\mathbf{B}^T$  and  $\mathbf{R}^T$  into dense matrices by merging non-overlapping sparse columns that do not contribute to the same nonzero in the output matrix. They show that this approach would incur less memory traffic than performing sparse inner products by a factor of  $n/n_{color}$  where  $n_{color}$  is the number of colors used for matrix coloring. However, they do not analyze the memory traffic of other formulations of SpGEMM, which are known to outperform sparse inner products [15]. In particular, a column-wise formulation of SpGEMM using CSC incurs only  $O(nnz/L + \text{flops})$  cache misses where  $L$  is the size of the cache line. Consider the matrix representing the Erdős-Rényi graph  $G(n, p)$ , where each edge (nonzero) in the graph (matrix) is present with probability  $p$  independently from each other. For  $p = d/n$  where  $d \ll n$ , in expectation  $nd$  nonzeros are uniformly distributed in an  $n$ -by- $n$  sparse matrix. In that case, SpGEMM does  $O(d^2n)$  cache misses compared to the  $O(dn n_{color})$  cache misses of the algorithm by McCourt et al. Hence, the column-wise approach not only bypasses the need for coloring, it also performs better for  $d \leq n_{color}$ , which is a common case. Furthermore, their method requires precomputing the nonzero structure of the output matrix, which is asymptotically as hard as computing SpGEMM without coloring in the first place.

There has been a flurry of activity in developing algorithms and implementations of SpGEMM for Graphics Processing Units (GPUs). Among those, the algorithm of

Gremse et al. [26] uses the row-wise formulation of SpGEMM. By contrast, Dalton et al. [18] uses the data-parallel ESC (expansion, sorting, and contraction) formulation, which is based on outer products. One downside of the ESC formulation is that expansion might create  $O(\text{flops})$  intermediate storage in the worst case, depending on the number of additions performed immediately in shared memory when possible, which might be asymptotically larger than the sizes of the inputs and outputs combined. The recent work of Liu and Vinter is currently the fastest implementation on GPUs and it also addresses heterogenous CPU-GPU processors [36].

In distributed memory, under many definitions of scalability, all known parallel SpGEMM algorithms are unscalable due to increased communication costs relative to arithmetic operations. For instance, there is no way to keep the parallel efficiency ( $PE$ ) fixed for any constant  $1 \geq PE > 0$  as we increase the number of processors [34]. Recently, two attempts have been made to model the communication costs of SpGEMM in a more fine grained manner. Akbudak and Aykanat [3] proposed the first hypergraph model for outer-product formulation of SpGEMM. Unfortunately, a symbolic SpGEMM computation has to be performed initially as the hypergraph model needs full access to the computational pattern that forms the output matrix. Ballard et al. [6] recently proposed hypergraph models for a class of SpGEMM algorithms more general than Akbudak and Aykanat considered. Their model also requires the sparsity structure of the output matrix and the number of vertices in the hypergraph is  $O(\text{flops})$ , making the approach impractical.

In terms of in-node parallelism via multithreading, there has been relatively little work. Gustavson’s algorithm is not thread scalable because its intermediate working set size is  $O(n)$  per thread, requiring a total of  $O(nt)$  intermediate storage, which can be larger than the matrices themselves for high thread counts. This intermediate data structure is called the sparse accumulator (SPA) [24]. Nevertheless, it is possible to get good performance out of a multithreaded parallelization of Gustavson’s algorithm in current platforms, provided that accesses to SPA are further “blocked” for matrices with large dimensions, in order to decrease cache miss rates. In a recent work, this is achieved by partitioning the data structure of the second matrix  $\mathbf{B}$  by columns [42].

We also mention that there has been significant research devoted to dense matrix multiplication in distributed-memory settings. In particular, the development of so-called 3D algorithms for dense matrix multiplication spans multiple decades; see [21, 30, 39, 44] and the references therein. Many aspects of our 3D algorithm for sparse matrix multiplication are derived from the dense case, though there are important differences as we detail below.

**4. Distributed-memory SpGEMM.** We categorize algorithms based on how they partition “work” (scalar multiplications) among processes, as we first advocated recently [5]. The work required by SpGEMM can be conceptualized by a cube that is sparsely populated by “voxels” that correspond to nonzero scalar multiplications. The algorithmic categorization is based on how these voxels are assigned to processes, which is illustrated in Figure 4.1. 1D algorithms assign a block of  $n$ -by- $n$ -by-1 “layers” of this cube to processes. In practice, this is realized by having each process store a block of rows or columns of an  $m$ -by- $n$  sparse matrix, though the 1D/2D/3D categorization is separate from the data distribution. With correctly chosen data distributions, 1D algorithms communicate entries of only one of the three matrices.

2D algorithms assign a set of 1-by-1-by- $n$  “fibers” of this cube to processes. In many practical realizations of 2D algorithms, processes are logically organized as a rectangular  $p = p_r \times p_c$  process grid, so that a typical process is named  $P(i, j)$ . Subma-

trices are assigned to processes according to a 2D block decomposition: For a matrix  $\mathbf{M} \in \mathbb{S}^{m \times n}$ , processor  $P(i, j)$  stores the submatrix  $\mathbf{M}_{ij}$  of dimensions  $(m/p_r) \times (n/p_c)$  in its local memory. With consistent data distributions, 2D algorithms communicate entries of two of the three matrices.

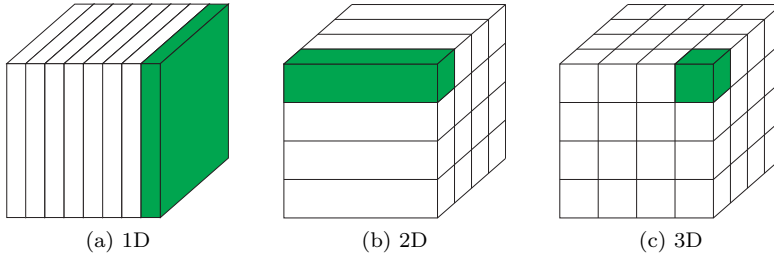


Fig. 4.1: Partitioning the work cube to processes. Image reproduced for clarity [5].

3D algorithms assign subcubes (with all 3 dimensions shorter than  $n$ ) to processes, which are typically organized on a  $p = p_r \times p_c \times p_l$  grid and indexed by  $P(i, j, k)$ . 3D algorithms communicate entries of  $\mathbf{A}$  and  $\mathbf{B}$ , as well as the (partial sums of the) intermediate products of  $\mathbf{C}$ . While many ways of assigning submatrices to processes on a 3D process grid exist, including replicating each  $\mathbf{A}_{ij}$  along the process fiber  $P(i, j, :)$ , our work focuses on a memory-friendly split decomposition. In this formulation,  $P(i, j, k)$  owns the following  $m/p_r \times n/(p_c p_l)$  submatrix of  $\mathbf{A} \in \mathbb{S}^{m \times n}$ :

$$\mathbf{A}(im/p_r : (i+1)m/p_r - 1, jn/p_c + kn/(p_c p_l) : jn/p_c + (k+1)n/(p_c p_l) - 1).$$

The distribution of matrices  $\mathbf{B}$  and  $\mathbf{C}$  are analogous. This distribution is memory friendly because it does not replicate input or output matrix entries, which is in contrast to many 3D algorithms where the input or the output is explicitly replicated.

**4.1. Sparse SUMMA Algorithm.** We briefly remind the reader of the Sparse SUMMA algorithm [11] for completeness as it will form the base of our 3D discussion. Sparse SUMMA is based on one formulation of the dense SUMMA algorithm [23]. The processes are logically organized on a  $p_r \times p_c$  process grid. The algorithm proceeds in stages where each stage involves the broadcasting of  $n/p_r \times b$  submatrices of  $\mathbf{A}$  by their owners along their process row, and the broadcasting of  $b \times n/p_c$  submatrices of  $\mathbf{B}$  by their owners along their process column. The recipients multiply the submatrices they received to perform a rank- $b$  update on their piece of the output matrix  $\mathbf{C}$ . The rank- $b$  update takes the form of a *merge* in the case of sparse matrices; several rank- $b$  updates can be done together using a multiway merge as described in Section 4.3. We will refer to this as a “SUMMA stage” for the rest of the paper. Here,  $b$  is a blocking parameter, which can be as large as the inner submatrix dimension. A more complete description of the algorithm and its general implementation for rectangular matrices and process grids can be found in an earlier work [14].

**4.2. In-node Multithreaded SpGEMM Algorithm.** Our previous work [12] shows that the standard compressed sparse column or row (CSC or CSR) data structures are too wasteful for storing the local submatrices arising from a 2D decomposition. This is because the local submatrices are *hypersparse*, meaning that the

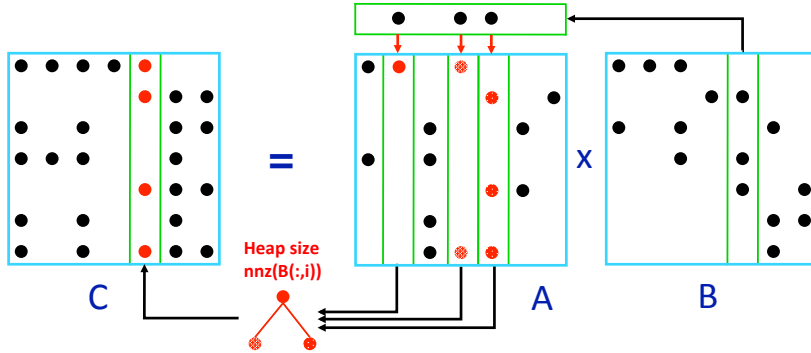


Fig. 4.2: Multiplication of sparse matrices stored by columns [12]. Columns of  $\mathbf{A}$  are accumulated as specified by the non-zero entries in a column of  $\mathbf{B}$  using a priority queue (heap) indexed by the row indices. The contents of the heap are stored into a column of  $\mathbf{C}$  once all required columns are accumulated.

ratio of nonzeros to dimension asymptotically approaches zero as the number of processors increase. The total memory across all processors for CSC format would be  $O(n\sqrt{p} + nnz)$ , as opposed to  $O(n + nnz)$  memory to store the whole matrix in CSC on a single processor.

This observation applies to 3D algorithms as well because their execution is reminiscent of running a 2D algorithm on each processor layer  $P(:, :, k)$ . Thus, local data structures used within 2D and 3D algorithms must respect hypersparsity.

Similarly, any algorithm whose complexity depends on matrix dimension, such as Gustavson's serial SpGEMM algorithm, is asymptotically too wasteful to be used as a computational kernel for multiplying the hypersparse submatrices. We use HEAP-SPGEMM, first presented as Algorithm 2 of our earlier work [12], which operates on the strictly  $O(nnz)$  doubly compressed sparse column (DCSC) data structure, and its time complexity does not depend on the matrix dimension. DCSC [12] is a further compressed version of CSC where repetitions in the column pointers array, which arise from empty columns, are not allowed. Only columns that have at least one nonzero are represented, together with their column indices. DCSC is essentially a sparse array of sparse columns, whereas CSC is a dense array of sparse columns. Although not part of the essential data structure, DCSC can support fast column indexing by building an AUX array that contains pointers to nonzero columns (columns that have at least one nonzero element) in linear time.

Our HEAPSPGEMM uses a heap-assisted column-by-column formulation whose time complexity is

$$\sum_{j=0}^{nzc(\mathbf{B})} O(\text{flops}(\mathbf{C}(:, j)) \log nnz(\mathbf{B}(:, j))),$$

where  $nzc(\mathbf{B})$  is the number of columns of  $\mathbf{B}$  that are not entirely zero,  $\text{flops}(\mathbf{C}(:, j))$  is the number of nonzero multiplications and additions required to generate the  $j$ th column of  $\mathbf{C}$ . The execution of this algorithm is illustrated in Figure 4.2, which differs from Gustavson's formulation in its use of a heap (priority queue) as opposed to a

sparse accumulator (SPA).

Our formulation is more suitable for multithreaded execution where we parallelize over the columns of  $\mathbf{C}$  and each thread computes  $\mathbf{A}$  times a subset of the columns of  $\mathbf{B}$ . SPA is an  $O(n)$  data structure, hence a multithreaded parallelization over columns of  $\mathbf{C}$  of the SPA-based algorithm would require  $O(nt)$  space for  $t$  threads. By contrast, since each heap in HEAPSPGEMM is of size  $O(nnz(\mathbf{B}(:,j)))$ , the total temporary memory requirements of our multithreaded algorithm are always strictly smaller than the space required to hold one of the inputs, namely  $\mathbf{B}$ .

**4.3. Multithreaded Multiway Merging and Reduction.** Each stage of Sparse SUMMA generates partial result matrices that are summed together at the end of all stages to obtain the final result  $\mathbf{C}$ . In the 3D algorithm discussed in Section 4.4, we also split  $\mathbf{C}$  across fibers of 3D grid, and the split submatrices are summed together by each process on the fiber. To efficiently perform these two summations, we represent the intermediate matrices as lists of triples, where each triple  $(i, j, val)$  stores the row index, column index, and value of a nonzero entry, respectively. Each list of triples is kept sorted lexicographically by the  $(j, i)$  pair so that the  $j$ th column comes before the  $(j+1)$ st column. We then perform the summation of sparse matrices by merging the lists of triples that represent the matrices. The merging also covers the reduction of triples with repeated indices.

To perform a  $k$ -way merge on  $k$  lists of triples  $T_1, T_2, \dots, T_k$ , we maintain a heap of size  $k$  that stores the current lexicographically minimum entry, based on  $(j, i)$  pairs, in each list of triples. In addition to keeping triples, the heap also stores the index of the source list from where each triple was inserted into the heap. The multiway merge routine finds the minimum triple  $(i^*, j^*, val^*)$  from the heap and merges it into the result. When the previous triple in the merged list has the same pair of indices  $(i^*, j^*)$ , the algorithm simply adds the values of these two triples, reducing the index-value pairs with repeated indices. If  $(i^*, j^*, val^*)$  is originated from  $T_l$ , the next triple from  $T_l$  is inserted into the heap. Hence, the time complexity of a  $k$ -way merge is

$$\sum_{l=1}^k O(nnz(T_l) \log k),$$

where  $nnz(T_l)$  is the number of nonzero entries in  $T_l$ . When multithreading is employed, each thread merges a subset of columns from  $k$  lists of triples using the same  $k$ -way merge procedure described earlier. If a thread is responsible for columns  $j_p$  to  $j_q$ , these column indices are identified from each list via a binary search. For better load balance, we ensure there is enough parallel slackness [46] by splitting the lists into more parts than the number of threads and merging the columns in each part by a single thread. In our experiments, we created  $4t$  parts when using  $t$  threads and employed dynamic thread scheduling.

**4.4. Split-3D-SpGEMM Algorithm.** Our parallel algorithm is an iterative 3D algorithm that splits the submatrices along the third process grid dimension (of length  $p_l$ ). This way, while there is no direct relationship between the size of the third process dimension and the extra memory required for *the input*, the extra memory required by *the output* is sparsity-structure dependent. If the output is sparse enough so that there are few or no intermediate products with repeated indices, then no extra memory is required. Recall that entries with repeated indices arise when more than one scalar multiplication  $a_{ik}b_{kj}$  that results in a nonzero value contributes to the same output element  $c_{ij}$ . The pseudocode of our algorithm, SPLIT-3D-SPGEMM, is shown



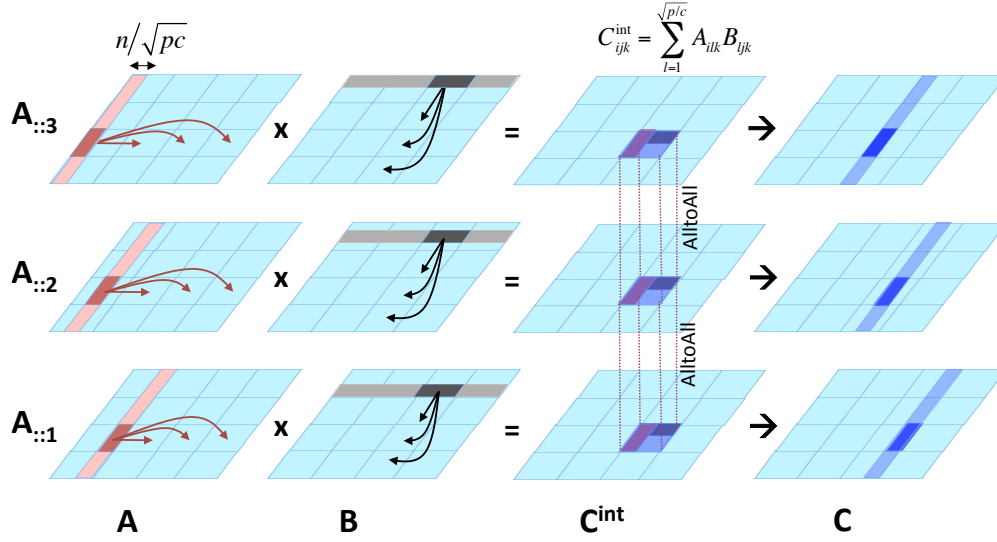


Fig. 4.3: Execution of the Split-3D-SpGEMM algorithm for sparse matrix-matrix multiplication  $\mathbf{C} = \mathbf{A} \cdot \mathbf{B}$  on a  $\sqrt{p/c} \times \sqrt{p/c} \times c$  process grid. Matrices  $\mathbf{A}$ ,  $\mathbf{B}$ , and  $\mathbf{C}^{int}$  matrices are shown during the first stage of the algorithm execution (the broadcast and the local update, i.e. one “SUMMA stage”). The transition from  $\mathbf{C}^{int}$  to  $\mathbf{C}$  happens via an all-to-all followed by a local merge, after all  $\sqrt{p/c}$  SUMMA stages are completed.

in Algorithm 2 for the simplified case of  $p_r = p_c = \sqrt{p/c}$  and  $p_l = c$ . The execution of the algorithm is illustrated in Figure 4.3.

---

**Algorithm 2** Operation  $\mathbf{C} \leftarrow \mathbf{A}\mathbf{B}$  using Split-3D-SpGEMM

---

**Input:**  $\mathbf{A} \in \mathbb{S}^{m \times l}$ ,  $\mathbf{B} \in \mathbb{S}^{l \times n}$ : matrices on a  $\sqrt{p/c} \times \sqrt{p/c} \times c$  process grid

**Output:**  $\mathbf{C} \in \mathbb{S}^{m \times n}$ : the product  $\mathbf{A}\mathbf{B}$ , similarly distributed.

- 1: **procedure** SPLIT-3D-SPGEMM( $\mathbf{A}, \mathbf{B}, \mathbf{C}$ )
  - 2:    $locinndim = l/\sqrt{pc}$  ▷ inner dimension of local submatrices
  - 3:   **for** all processes  $P(i, j, k)$  **in parallel do**
  - 4:      $\hat{\mathbf{B}}_{ijk} \leftarrow \text{ALLTOALL}(\mathbf{B}_{ij}, P(i, j, :))$  ▷ redistribution of  $\mathbf{B}$  across layers
  - 5:     **for**  $r = 1$  to  $\sqrt{p/c}$  **do** ▷  $r$  is the broadcasting process column and row
  - 6:       **for**  $q = 1$  to  $locinndim/b$  **do** ▷  $b$  evenly divides  $locinndim$
  - 7:          $locindices = (q-1)b : qb - 1$
  - 8:          $\mathbf{A}^{rem} \leftarrow \text{BROADCAST}(\mathbf{A}_{irk}(:, locindices), P(i, :, k))$
  - 9:          $\mathbf{B}^{rem} \leftarrow \text{BROADCAST}(\hat{\mathbf{B}}_{rjk}(locindices, :), P(:, j, k))$
  - 10:         $\mathbf{C}_{ij:}^{int} \leftarrow \mathbf{C}_{ij:}^{int} + \text{HEAPSPGEMM}(\mathbf{A}^{rem}, \mathbf{B}^{rem})$
  - 11:      $\mathbf{C}_{ijk}^{int} \leftarrow \text{ALLTOALL}(\mathbf{C}_{ij:}^{int}, P(i, j, :))$
  - 12:      $\mathbf{C}_{ijk} \leftarrow \text{LOCALMERGE}(\mathbf{C}_{ijk}^{int})$
- 

The  $\text{BROADCAST}(\mathbf{A}_{irk}, P(i, :, k))$  syntax means that the owner of  $\mathbf{A}_{irk}$  becomes the root and broadcasts its submatrix to all the processes on the  $i$ th process row of

the  $k$ th process layer. Similarly for  $\text{BROADCAST}(\hat{\mathbf{B}}_{rjk}, P(:, j, k))$ , the owner of  $\hat{\mathbf{B}}_{rjk}$  broadcasts its submatrix to all the processes on the  $j$ th process column of the  $k$ th process layer. In line 7, we find the local column (for  $\mathbf{A}$ ) and row (for  $\hat{\mathbf{B}}$ ) ranges for matrices that are to be broadcast during that iteration. They are significant only at the broadcasting processes, which can be determined implicitly from the first parameter of  $\text{BROADCAST}$ . In practice, we index  $\hat{\mathbf{B}}$  by columns as opposed to rows in order to obtain the best performance from the column-based DCSC data structure. To achieve this,  $\hat{\mathbf{B}}$  gets locally transposed during redistribution in line 4. Using DCSC, the expected cost of fetching  $b$  consecutive columns of a matrix  $\mathbf{A}$  is  $b$  plus the size (number of nonzeros) of the output. Therefore, the algorithm asymptotically has the same computation cost for all blocking parameters  $b$ .

$\mathbf{C}_{ij}^{int}$  is the intermediate submatrix that contains nonzeros that can potentially belong to all the processes on the  $(i, j)$ th fiber  $P(i, j, :)$ . The ALLTOALL call in line 11 packs those nonzeros and sends them to their corresponding owners in the  $(i, j)$ th fiber. This results in  $\mathbf{C}_{ijk}^{int}$  for each process  $P(i, j, k)$ , which contains only the nonzeros that belong to that process.  $\mathbf{C}_{ijk}^{int}$  possibly contains repeated indices (i.e. multiple entries with the same index) that need to be merged and summed by the LOCALMERGE call in line 12, resulting in the final output.

In contrast to dense matrix algorithms [21, 30, 39, 44], our sparse 3D formulation requires a more lenient trade-off between bandwidth-related communication costs and memory requirements. As opposed to increasing the storage requirements by a factor of  $p_l$ , the relative cost of the 3D formulation is  $nnz(\mathbf{C}^{int})/nnz(\mathbf{C})$ , which is always upper bounded by  $\text{flops}(\mathbf{A}, \mathbf{B})/nnz(\mathbf{C})$ .

**4.5. Communication Analysis of the Split-3D-SpGEMM Algorithm.**

For our complexity analysis, the previous work [5] assumed that nonzeros of sparse  $n$ -by- $n$  input matrices are independently and identically distributed, with  $d > 0$  nonzeros per row and column on the average. The sparsity parameter  $d$  simplifies analysis by making different terms in the complexity comparable to each other. However, in order to capture the performance of more general matrix-matrix multiplication, we will analyze parallel complexity directly in terms of flops and the number of nonzeros in  $\mathbf{A}$ ,  $\mathbf{B}$ , and  $\mathbf{C}$  without resorting to the sparsity parameter  $d$ .

Our algorithm can run on a wide range of configurations on a virtual 3D  $p = p_r \times p_c \times p_l$  process grid. To simplify the analysis, we again assume that each 2D layer of the 3D grid is square, i.e.  $p_r = p_c$  and we use  $c$  to denote the third dimension. Thus, we assume a  $\sqrt{p/c} \times \sqrt{p/c} \times c$  process grid.

The communication in Algorithm 2 consists of collective operations being performed on disjoint process fibers: simultaneous broadcasts in the first two process grid dimensions at line 8 and line 9 and simultaneous all-to-alls in the third process grid dimension at line 4 and line 11. We use the notation  $T_{\text{bcast}}(w, \hat{p}, \nu, \mu)$  and  $T_{\text{a2a}}(w, \hat{p}, \nu, \mu)$  to denote the costs of broadcast and all-to-all, where  $w$  is the size of the data (per processor) in matrix elements,  $\hat{p}$  is the number of processes participating in the collective,  $\nu$  is the number of simultaneous collectives, and  $\mu$  is the number of processes per node. Parameters  $\nu$  and  $\mu$  capture resource contention: the number of simultaneous collectives affects contention for network bandwidth, and the number of processes per node affects contention for the network interface card on each node.

In general, these cost functions can be approximated via microbenchmarks for a given machine and MPI implementation, though they can vary over different node allocations as well. If we ignore resource contention, with  $\nu = 1$  and  $\mu = 1$ , then the

costs are typically modeled [17] as

$$T_{\text{broadcast}}(w, \hat{p}, 1, 1) = \alpha \cdot \log \hat{p} + \beta \cdot w \frac{\hat{p} - 1}{\hat{p}}$$

and

$$T_{\text{a2a}}(w, \hat{p}, 1, 1) = \alpha \cdot (\hat{p} - 1) + \beta \cdot w \frac{\hat{p} - 1}{\hat{p}}.$$

The all-to-all cost assumes a point-to-point algorithm, minimizing bandwidth cost at the expense of higher latency cost; see [5, Section 2.2] for more details on the tradeoffs within all-to-all algorithms.

The time spent in communication is then given by

$$\begin{aligned} & T_{\text{a2a}} \left( \frac{\text{nnz}(\mathbf{B})}{p}, c, \frac{p}{c}, \mu \right) + \\ & \frac{n}{bc} \cdot T_{\text{broadcast}} \left( \frac{b}{n} \cdot \frac{\text{nnz}(\mathbf{A})}{\sqrt{p/c}}, \sqrt{p/c}, \sqrt{pc}, \mu \right) + \\ & \frac{n}{bc} \cdot T_{\text{broadcast}} \left( \frac{b}{n} \cdot \frac{\text{nnz}(\mathbf{B})}{\sqrt{p/c}}, \sqrt{p/c}, \sqrt{pc}, \mu \right) + \\ & T_{\text{a2a}} \left( \frac{\text{flops}(\mathbf{A}, \mathbf{B})}{p}, c, \frac{p}{c}, \mu \right). \end{aligned}$$

The amount of data communicated in the first three terms is the average over all processes and is accurate only if the nonzeros of the input matrices are evenly distributed across all blocks. The amount of data communicated in the last term is an upper bound on the average; the number of output matrix entries communicated by each process is likely less than the number of flops performed by that process (due to the reduction of locally repeated indices prior to communication). A lower bound for the last term is given by replacing  $\text{flops}(\mathbf{A}, \mathbf{B})$  with  $\text{nnz}(\mathbf{C})$ .

If we ignore resource contention, the communication cost is

$$\alpha \cdot O \left( \frac{n}{bc} \log(p/c) + c \right) + \beta \cdot O \left( \frac{\text{nnz}(\mathbf{A}) + \text{nnz}(\mathbf{B})}{\sqrt{pc}} + \frac{\text{flops}(\mathbf{A}, \mathbf{B})}{p} \right),$$

where we have assumed that  $\text{nnz}(\mathbf{B}) \leq \text{flops}(\mathbf{A}, \mathbf{B})$ . Note that this expression matches the costs for Erdős-Rényi matrices, up to the choice of all-to-all algorithm [5], where  $\text{nnz}(\mathbf{A}) \approx \text{nnz}(\mathbf{B}) \approx dn$  and  $\text{flops}(\mathbf{A}, \mathbf{B}) \approx d^2n$ .

We make several observations based on this analysis of the communication. First, increasing  $c$  (the number of layers) results in less time spent in broadcast collectives and more time spent in all-to-all collectives (note that if  $c = 1$  then no communication occurs in all-to-all). Second, increasing  $b$  (the blocking parameter) results in fewer collective calls but the same amount of data communicated; thus,  $b$  navigates a tradeoff between latency cost and local memory requirements (as well as greater possibility to overlap local computation and communication). Third, for a fixed number of cores, lower  $\mu$  (higher value of  $t$ ) will decrease network interface card contention and therefore decrease communication time overall.

	<b>Cray XK7 (Titan)</b>	<b>Cray XC30 (Edison)</b>
<b>Core</b>	<b>AMD Interlagos</b>	<b>Intel Ivy Bridge</b>
Clock (GHz)	2.2	2.4
Private Cache (KB)	16+2048	32+256
DP GFlop/s/core	8.8	19.2
<b>Socket Arch.</b>	<b>Opteron 6172</b>	<b>Xeon E5-2695 v2</b>
Cores per socket	16	12
Threads per socket	16	24 <sup>1</sup>
L3 cache per socket	2×8 MB	30 MB
<b>Node Arch.</b>	<b>Hypertransport</b>	<b>QPI (8 GT/s)</b>
Sockets/node	1	2
STREAM BW <sup>2</sup>	31 GB/s	104 GB/s
Memory per node	32 GB	64 GB
<b>Interconnect</b>	<b>Gemini (3D Torus)</b>	<b>Aries (Dragonfly)</b>

Table 5.1: Overview of Evaluated Platforms. <sup>1</sup>Only 12 threads were used. <sup>2</sup>Memory bandwidth is measured using the STREAM copy benchmark per node.

**5. Experimental Results.** We evaluate our algorithms on two supercomputers: Cray XC30 at NERSC (Edison) [22], and Cray XK6 at ORNL (Titan) [45]. Architectural details of these computers are listed in Table 5.1. In our experiments, we ran only on the CPUs and did not utilize Titan’s GPU accelerators.

In both supercomputers, we used Cray’s MPI implementation, which is based on MPICH2. Both chip architectures achieve memory parallelism via hardware prefetching. On Titan, we compiled our code using GCC C++ compiler version 4.6.2 with `-O2 -fopenmp` flags. On Edison, we compiled our code using the Intel C++ compiler (version 14.0.2) with the options `-O2 -no-ipo -openmp`. In order to ensure better memory affinity to NUMA nodes of Edison and Titan, we used `-cc depth` or `-cc numa_node` options when submitting jobs. For example, to run the 3D algorithm on a  $8 \times 8 \times 4$  process grid with 6 threads, we use the following options on Edison: `aprun -n 256 -d 6 -N 4 -S 2 -cc depth`. In our experiments, we always allocate cores needed for a particular configuration of 3D algorithms, i.e., to run the 3D algorithm on  $\sqrt{p/c} \times \sqrt{p/c} \times c$  process grid with  $t$  threads per process, we allocate  $pt$  cores and run  $p$  MPI processes on the allocated cores.

Several software libraries support SpGEMM. For GPUs, CUSP and CUSparse implement SpGEMM. For shared-memory nodes, MKL implements SpGEMM. Trilinos package implements distributed memory SpGEMM [29], which uses a 1D decomposition for its sparse matrices. In this paper, we compared the performance of 2D and 3D algorithms with SpGEMM in Cray-Trilinos package (version 11.6.1.0) available in NERSC computers, which features significant performance improvements over earlier versions. Sparse SUMMA is the 2D algorithm that had been published before [11] without in-node multithreading, and Split-3D-SpGEMM is the 3D algorithm first presented here. Sometimes we will drop the long names and just use 2D and 3D for abbreviation.

In our experiments, we used both synthetically generated matrices, as well as real matrices from several different sources. In Section 5.2, we benchmark square matrix multiplication. We use R-MAT [16], the Recursive MATrix generator to generate three different classes of synthetic matrices: (a) G500 matrices representing graphs

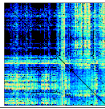
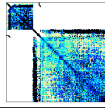
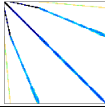
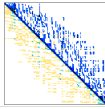
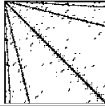
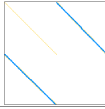
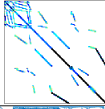
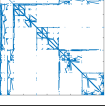
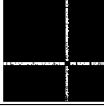
Name		Dimensions	nnz/row
Description	Spy Plot	Nonzeros	symmetric
<b>mouse_gene</b> Gene network		45K×45K 28.9M	642 ✓
<b>ldoor</b> structural problem		952K×952K 46.5M	48.8 ✓
<b>dielFilterV3real</b> electromagnetics problem		1.1M×1.1M 89.3M	81.2 ✓
<b>cage15</b> DNA electrophoresis		5.15M×5.15M 99.2M	19.3
<b>delaunay_n24</b> Delaunay triangulation		16.77M×16.77M 100.6M	6 ✓
<b>nlpkkt160</b> indefinite KKT matrix		8.34M×8.34M 229.5M	27.5 ✓
<b>HV15R</b> 3D engine fan		2.01M×2.01M 283M	141.5
<b>NaluR3</b> Low Mach fluid flow		17.6M×17.6M 474M	26.9 ✓
<b>it-2004</b> web crawl of .it domain		41.29M×41.29M 1,150M	27.8

Fig. 5.1: Structural information on the sparse matrices used in our experiments. All matrices are from the University of Florida sparse matrix collection [20], except NaluR3, which is a matrix from low Mach number, turbulent reacting flow problem [35]. For the Florida matrices, we consider the explicit zero entries to be nonzeros and update the nnz of the matrices accordingly.

with skewed degree distributions from Graph 500 benchmark [1], (b) SSCA matrices from HPCS Scalable Synthetic Compact Applications graph analysis (SSCA#2) benchmark [2], and (c) ER matrices representing Erdős-Rényi random graphs. We use the following R-MAT seed parameters to generate these matrices: (a)  $a = .57$ ,  $b = c = .19$ , and  $d = .05$  for G500, (b)  $a = .6$ , and  $b = c = d = .4/3$  for SSCA, and (c)  $a = b = c = d = .25$  for ER. A scale  $n$  synthetic matrix is  $2^n$ -by- $2^n$ . On average, G500 and ER matrices have 16 nonzeros, and SSCA matrices have 8 nonzeros per row

Table 5.2: Statistics about squaring real matrices and multiplying each matrix with its restriction operator  $\mathbf{R}$ .

Matrix ( $\mathbf{A}$ )	$nnz(\mathbf{A})$	$nnz(\mathbf{A}^2)$	$nnz(\mathbf{R})$	$nnz(\mathbf{R}^T \mathbf{A})$	$nnz(\mathbf{R}^T \mathbf{A} \mathbf{R})$
<b>mouse_gene</b>	28,967,291	482,594,045	45,101	2,904,560	402,200
<b>ldoor</b>	46,522,475	145,422,935	952,203	2,308,794	118,093
<b>dielFilterV3real</b>	89,306,020	688,649,400	1,102,824	4,316,781	100,126
<b>cage15</b>	99,199,551	929,023,247	5,154,859	46,979,396	17,362,065
<b>delaunay_n24</b>	100,663,202	347,322,258	16,777,216	41,188,184	15,813,983
<b>nlpkt160</b>	229,518,112	1,241,294,184	8,345,600	45,153,930	3,645,423
<b>HV15R</b>	283,073,458	1,768,066,720	2,017,169	10,257,519	1,400,666
<b>NaluR3</b>	473,712,505	2,187,662,967	17,598,889	77,245,697	7,415,297
<b>it-2004</b>	1,150,725,436	14,045,664,641	41,291,594	89,870,859	26,847,490

and column. We applied a random symmetric permutation to the input matrices to balance the memory and the computational load. In other words, instead of storing and computing  $\mathbf{C} = \mathbf{A}\mathbf{B}$ , we compute  $\mathbf{PCP}^T = (\mathbf{PAP}^T)(\mathbf{PBP}^T)$ . All of our experiments are performed on double-precision floating-point inputs, and matrix indices are stored as 64-bit integers.

In Section 5.3, we benchmark the matrix multiplication corresponding to the restriction operation that is used in AMG. Since AMG on graphs coming from physical problems is an important case, we include several matrices from the Florida Sparse Matrix collection [20] in our experimental analysis. In addition, since AMG restriction is computationally isomorphic to the graph contraction operation performed by multilevel graph partitioners [28], we include a few matrices representing real-world graphs.

The characteristics of the real test matrices are shown in Table 5.1. Statistics about squaring real matrices and multiplying each matrix with its restriction operator  $\mathbf{R}$  is given in Table 5.2.

**5.1. Intra-node Performance.** Our 3D algorithm exploits intra-node parallelism in two computationally intensive functions: (a) local HEAPSPGEMM performed by each MPI process at every SUMMA stage, and (b) multiway merge performed at the end of all SUMMA stages. As mentioned before, HEAPSPGEMM returns a set of intermediate triples that are kept in memory and merged at the end of all SUMMA stages. In this section, we only show the intra-node scalability of these two functions and compare them against an MKL and a GNU routine.

**5.1.1. Multithreaded HEAPSPGEMM Performance.** We study intra-node scalability of local SpGEMM by running a single MPI process on a socket of a node and varying the number of threads from one to the maximum number of threads available in a socket. We compare the performance of HEAPSPGEMM with MKL routine `mk1_csrmultcsr`. To expedite the multiway merge that is called on the output of HEAPSPGEMM, we always keep column indices sorted in increasing order within each row. Hence, we ask `mk1_csrmultcsr` to return sorted output. We show the performance of HEAPSPGEMM and MKL in Figure 5.2 where these functions are used to multiply (a) two randomly generated scale 16 G500 matrices, and (b) Cage12 matrix by itself. On 12 threads of Edison, HEAPSPGEMM achieves  $8.5\times$  speedup for scale 16 G500 and  $8.7\times$  speedup for Cage12, whereas MKL achieves  $7.1\times$  speedup for scale 16 G500 and  $9.2\times$  speed for Cage12. Hence, HEAPSPGEMM scales as well

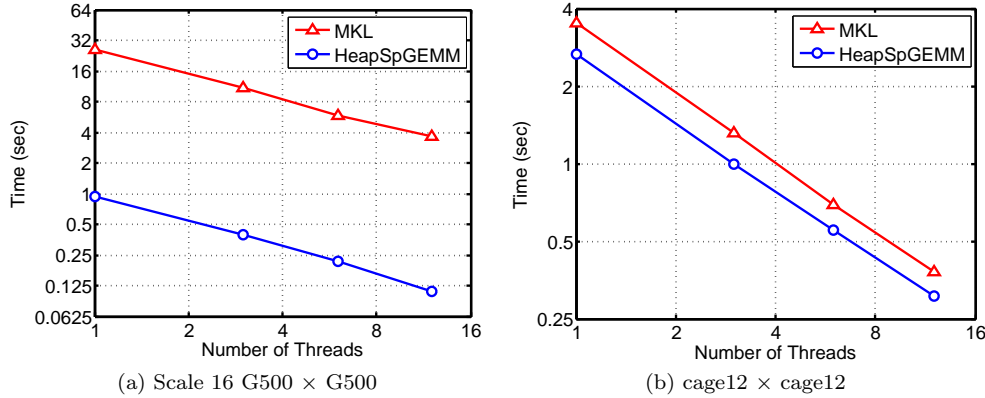


Fig. 5.2: Thread scaling of our HEAPSPGEMM and the MKL routine `mk1_csrmultcsr` on 1,3,6,12 threads (with column indices sorted in the increasing order for each row) when squaring a matrix on a single socket of Edison.

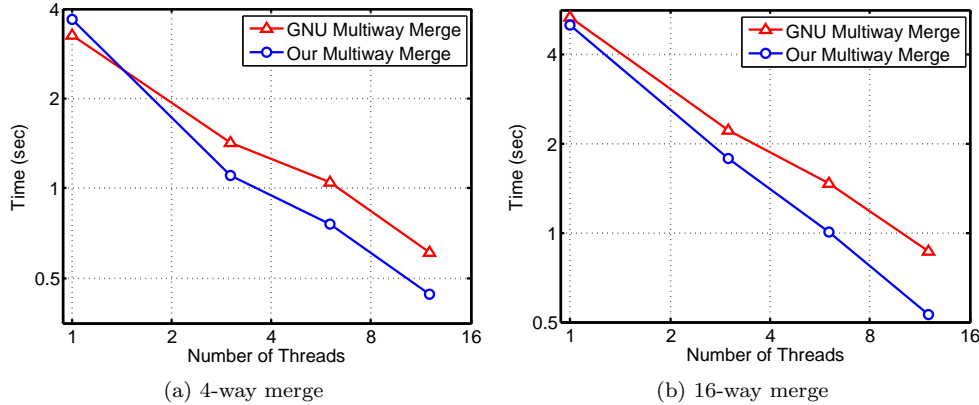


Fig. 5.3: Thread scaling of our multiway merge and GNU multiway merge routine `_gnu_parallel::multiway_merge` augmented with a procedure that reduces repeated indices: (a) in squaring scale 21 G500 matrix on  $4 \times 4$  process grid and (b) in squaring scale 26 G500 matrix on  $16 \times 16$  process grid on Edison.

as MKL. However, for these matrices, HEAPSPGEMM runs faster than MKL on any concurrency with up to 33-fold performance improvement for G500 matrices.

**5.1.2. Multithreaded Multiway Merge Performance.** On a  $\sqrt{p/c} \times \sqrt{p/c} \times c$  process grid, each MPI process performs two multiway-merge operations. The first one merges  $\sqrt{p/c}$  lists of triples computed in  $\sqrt{p/c}$  stages of SUMMA and the second one merges  $c$  lists of triples after splitting the previously merged list across layers. Since both merges are performed by the same function, we experimented the intra-node performance of multiway merge on a single layer ( $c=1$ ). For this experiment, we allocate  $12p$  cores on Edison and run SUMMA on a  $\sqrt{p} \times \sqrt{p}$  process grid. Each MPI process is run on a socket using up to 12 available threads. Figure 5.3 shows the merge time needed by MPI rank 0 for a 4-way merge and a 16-way merge when

multiplying two G500 matrices on a  $4 \times 4$  and a  $16 \times 16$  grid, respectively. We compare the performance of our multiway merge with a GNU multiway merge routine `__gnu_parallel::multiway_merge`. However, the latter merge routine simply merges lists of triples keeping them sorted by column and row indices, but does not reduce triples with the same (row, column) pair. Hence, we reduce the repeated indices returned by `__gnu_parallel::multiway_merge` by a multithreaded reduction function and report the total runtime. From Figure 5.3 we observe that our routine performs both 4-way and 16-way merges faster than augmented GNU multiway merge for G500 matrices. On 12 threads of Edison, our multiway merge attains  $8.3\times$  speedup for 4-way merge and  $9.5\times$  speedup for 16-way merge. By contrast, the augmented GNU merge attains  $5.4\times$  and  $6.2\times$  speedups for 4-way and 16-way merges, respectively. We observe similar performances for other matrices as well.

**5.2. Square Sparse Matrix Multiplication.** In the first set of experiments, we square real matrices from Table 5.1 and multiply two structurally similar randomly generated matrices. This square multiplication is representative of the expansion operation used in the Markov clustering algorithm [47]. We explore an extensive set of parameters of Sparse SUMMA (2D) and Split-3D-SpGEMM (which is the main focus of this work), identify optimum parameters on different levels of concurrency, empirically explain where Split-3D-SpGEMM gains performance, and then show the scalability of Split-3D-SpGEMM for a comprehensive set of matrices.

**5.2.1. Performance of Different Variants of 2D and 3D Algorithms.** At first, we investigate the impact of multithreading and 3D algorithm on the performance of SpGEMM. For this purpose, we fix the number of cores  $p$  and multiply two sparse matrices with different combinations of thread counts  $t$  and number of layers  $c$ . Figure 5.4 shows strong scaling of squaring of `nlpkkt160` matrix on Edison. On lower concurrency ( $< 256$  cores), multithreading improves the performance of 2D algorithm, e.g., about  $1.5\times$  performance improvement with 6 threads on 256 cores in Figure 5.4. However, there is little or no benefit in using a 3D algorithm over a multithreaded 2D algorithm on lower concurrency because the processor grid in a layer becomes too small for 3D algorithms.

For better resolution on higher concurrency, we have not shown the runtime of 2D algorithms before 64 cores in Figure 5.4. For the completeness of our discussion, we briefly discuss the performance of our algorithm on lower concurrency and compare them against MKL and Matlab. Matlab uses an efficient CSC implementation of Gustavson’s algorithm. 2D non-threaded algorithm takes about 800 seconds on a single core and attains about  $50\times$  speedup when we go from 1 core to 256 cores, and 2D algorithm with 6 threads attains about  $25\times$  speedup when we go from 6 cores to 216 cores. By contrast, on a single core, MKL and Matlab take about 500 and 830 seconds, respectively to square randomly permuted `nlpkkt160` matrix (in Matlab, we keep explicit zero entries<sup>1</sup> to obtain the same number of nonzeros shown in Table 5.2). Therefore, the serial performance of Sparse SUMMA (2D) is comparable to that of MKL and Matlab. The best single node performance is obtained by multithreaded SpGEMM. Using 24 threads on 24 cores of a single node of Edison, MKL and HeapSPGEMM take about 32 and 30 seconds, respectively. We note that the above performance numbers depend significantly on nonzero structures of the input matrices. Here, we select `nlpkkt160` matrix for discussion because the number

---

<sup>1</sup>The default Matlab behavior is to remove entries with zero values when constructing a matrix using its `sparse(i,j,v)`



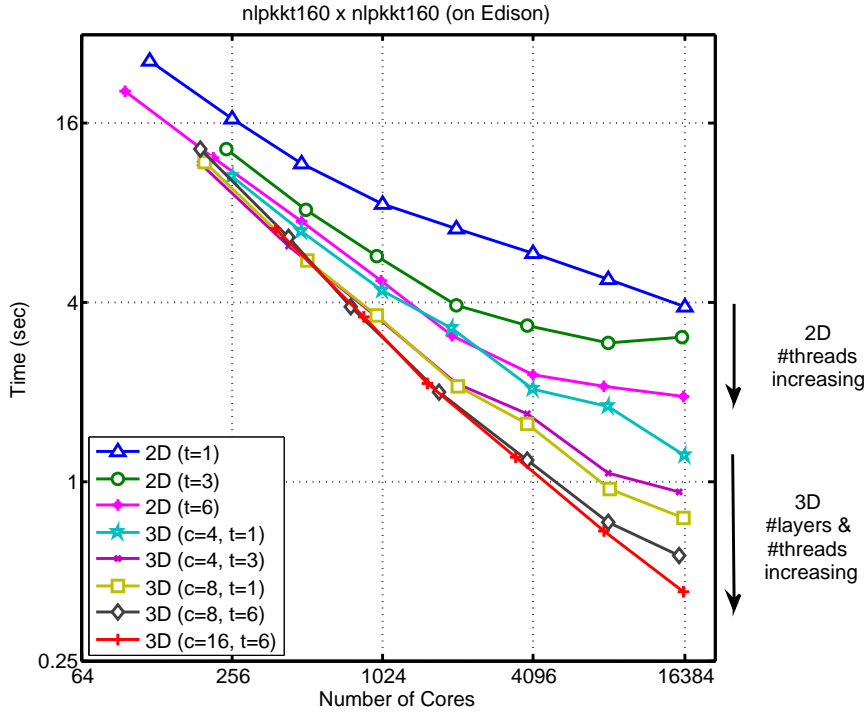


Fig. 5.4: Strong scaling of different variants of 2D (Sparse SUMMA) and 3D (Split-3D-SpGEMM) algorithms when squaring of `nlpkkt160` matrix on Edison. Performance benefits of the 3D algorithm and multithreading can be realized on higher concurrency. 2D non-threaded algorithm attains about  $50\times$  speedup when we go from 1 core to 256 cores, and 2D algorithm with 6 threads attains about  $25\times$  speedup when we go from 6 cores to 216 cores (not shown in the figure).

of nonzero in the square of `nlpkkt160` is about 1.2 billion (c.f. Table 5.2), requiring about 28GB of memory to store the result, which is close to the available single node memory of Edison.

The performance benefits of the 3D algorithm and multithreading become more dominant on higher concurrency. In Figure 5.4, when we increase  $p$  from 256 to 16,384 ( $64\times$  increase), non-threaded 2D and 3D ( $c=16, t=6$ ) algorithms run  $4\times$  and  $22\times$  faster, respectively. Consequently, on 16,384 cores, Split-3D-SpGEMM with  $c=16, t=6$  multiplies `nlpkkt160` matrix  $8\times$  faster than non-threaded 2D algorithm. We observe similar trend for other real and randomly generated matrices as well. For example, Split-3D-SpGEMM with  $c=16, t=6$  runs  $10\times$  faster than the Sparse SUMMA (2D) algorithm when squaring of `NaluR3` on 32,764 cores of Edison (Figure 5.5), and Split-3D-SpGEMM with  $c=16, t=8$  runs  $9.5\times$  faster than 2D algorithm when multiplying two scale 26 RMat matrices on 65,536 cores of Titan (Figure 5.6).

In fact, on higher concurrency, the time Split-3D-SpGEMM takes to multiply two square matrices decreases gradually with the increase of  $c$  and  $t$  as indicated on the right side of Figure 5.4. This trend is also observed in Figures 5.5 and 5.6. Therefore, we expect that using more threads and layers will be beneficial to gain performance on even higher concurrency.

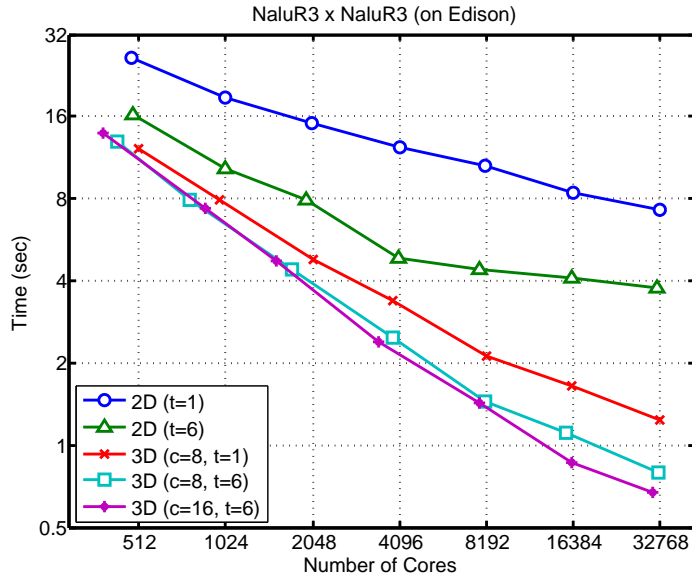


Fig. 5.5: Strong scaling of different variants of 2D and 3D algorithms when squaring of NaluR3 matrix on Edison. 3D algorithms are an order of magnitude faster than 2D algorithms on higher concurrency.

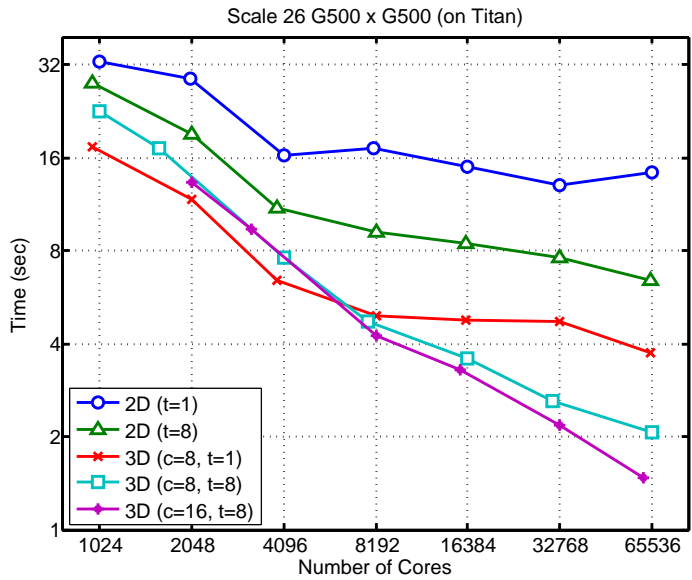


Fig. 5.6: Strong scaling of different variants of 2D and 3D algorithms on Titan when multiplying two scale 26 G500 matrices. 3D algorithm and multithreading improve performance of SpGEMM on higher concurrency

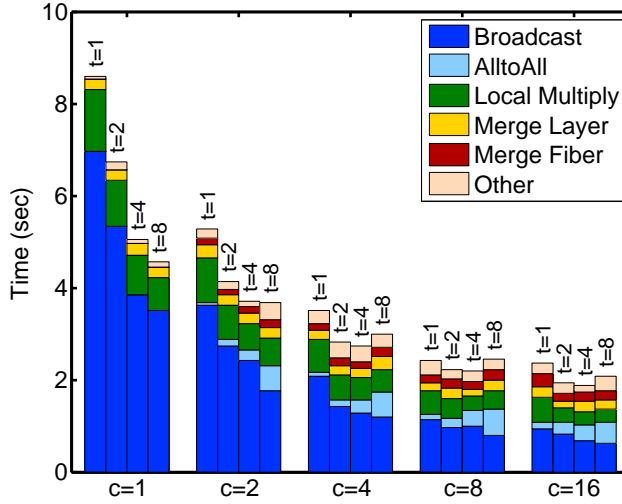


Fig. 5.7: Breakdown of runtime spent by Split-3D-SpGEMM for various  $(c, t)$  configurations on 8,192 cores of Titan when multiplying two scale 26 G500 matrices. The broadcast time (the most dominating term on high concurrency) decreases gradually with the increase of both  $c$  and  $t$ , which is the primary catalyst behind the improved performance of multithreaded 3D algorithms.

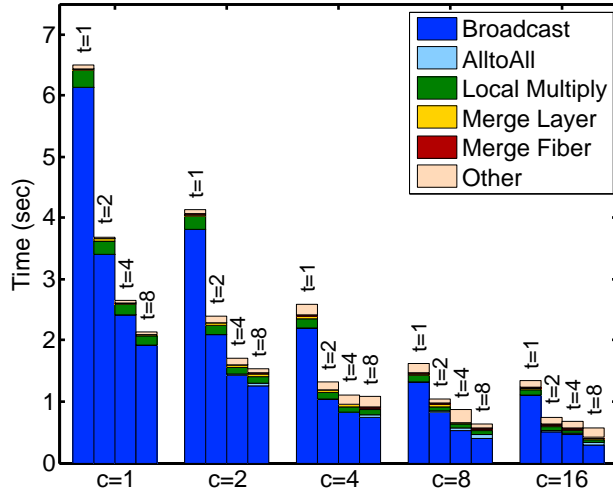


Fig. 5.8: Breakdown of runtime spent by Split-3D-SpGEMM for various  $(c, t)$  configurations on 32,768 cores of Titan when multiplying two scale 26 G500 matrices.

**5.2.2. Breakdown of Runtime.** To understand the performance of Split-3D-SpGEMM, we break down the time spent in communication and computation when multiplying two G500 graphs of scale 26 and show them in Figure 5.7 for 8,192 cores and Figure 5.8 for 32,768 cores on Titan. Here, “Broadcast” refers to the time needed to broadcast pieces of  $\mathbf{A}$  and  $\mathbf{B}$  within each layer, “AlltoAll” refers to the communication time needed to communicate pieces of  $\mathbf{C}$  across layers, “Local Multiply” is the

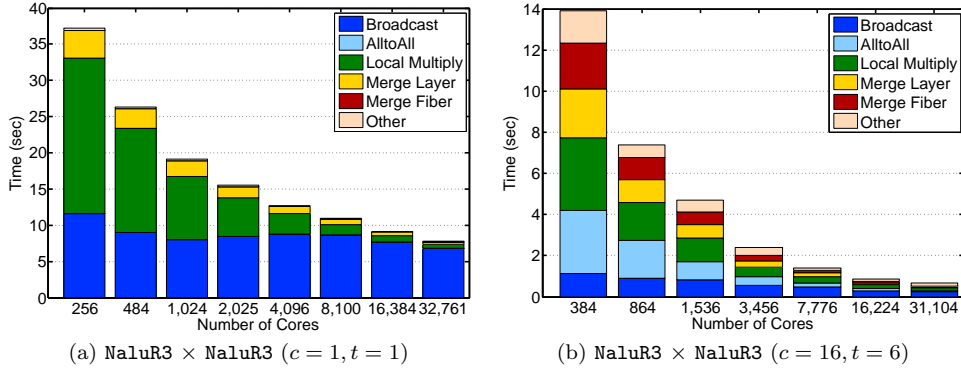


Fig. 5.9: Breakdown of runtime spent by (a) Sparse SUMMA (2D) algorithm with  $c = 1, t = 1$  and (b) Split-3D-SpGEMM algorithm with  $c = 16, t = 6$  to square `NaLuR3` on Edison.

time needed by multithreaded `HEAPSPGEMM`, “Merge Layer” is the time to merge  $\sqrt{p/c}$  lists of triples computed in  $\sqrt{p/c}$  stages of SUMMA within a layer, and “Merge Fiber” is the time to merge  $c$  lists of triples after splitting pieces of  $\mathbf{C}$  across processor fibers. For a fixed number of cores, the broadcast time gradually decreases with the increase of both  $c$  and  $t$ , because as we increase  $c$  and/or  $t$ , the number of MPI processes participating in broadcast within each process layer decreases. For example, in Figure 5.8, the broadcast time decreases by more than  $5\times$  from the leftmost bar to the rightmost bar. Since broadcast is the dominating term on higher concurrency, reducing it improves the overall performance of SpGEMM. However, for a fixed number of cores, the All2All time increases with  $c$  due to the increased processor count on the fiber. The All2All time also increases with  $t$  because each MPI process owns a bigger portion of the data, increasing the All2All communication cost per process. Therefore, increased All2All time might nullify the advantage of reduced broadcast time when we increase  $c$  and  $t$ , especially on lower concurrency. For example, using  $c > 4$  does not reduce the total communication time on 8,192 as shown in Figure 5.7.

Figure 5.7 and Figure 5.8 reveal that shorter communication time needed by Split-3D-SpGEMM makes it faster than Sparse SUMMA (2D) on higher concurrency. Figure 5.9(b) demonstrates that both communication and computation time scale well for Split-3D-SpGEMM with  $c = 16, t = 6$  when squaring `NaLuR3` on Edison. By contrast, communication time does not scale well for Sparse SUMMA (Figure 5.9(a)), which eventually limits the scalability of 2D algorithms on higher concurrency.

**5.2.3. Strong Scaling of Split-3D-SpGEMM.** In this subsection, we show the strong scaling of Split-3D-SpGEMM with the best parameters on higher concurrency ( $c = 16, t = 6$  on Edison, and  $c = 16, t = 8$  on Titan) when multiplying real and random matrices. Figure 5.10 shows the strong scaling of Split-3D-SpGEMM when squaring seven real matrices on Edison. When we go from 512 to 32,768 cores ( $64\times$  increase of cores), the average speedup of all matrices in Table 5.1 is about  $27\times$  (min:  $9\times$  for `deLaunay_n24`, max:  $52\times$  for `mouse_gene`, standard deviation: 16). We observe that Split-3D-SpGEMM scales better when multiplying larger (e.g., `it-2004`) and denser matrices (e.g., `mouse_gene` and `HV15R`) because of the availability of more work. By contrast, `deLaunay_n24` is the sparsest matrix in Table 5.1 with 6 nonzeros

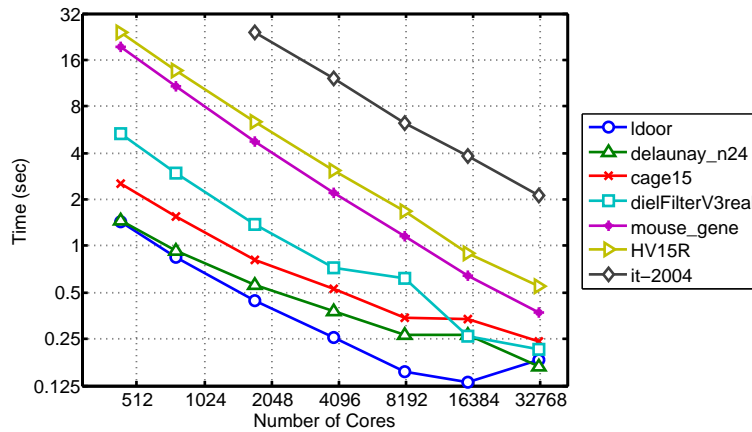


Fig. 5.10: Strong scaling of Split-3D-SpGEMM with  $c = 16, t = 6$  when squaring real matrices on Edison. Large (e.g., `it-2004`) and dense (e.g., `mouse_gene` and `HV15R`) matrices scale better than small and sparse (e.g., `delaunay_n24`) matrices.

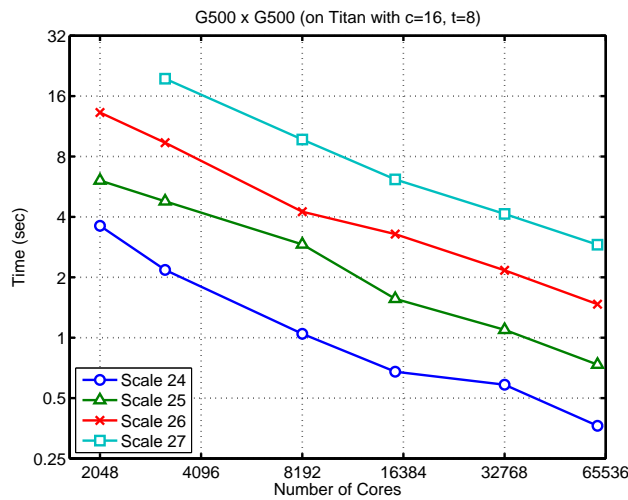


Fig. 5.11: Strong scaling of Split-3D-SpGEMM with  $c = 16, t = 8$  on Titan when multiplying two G500 matrices.

per column, and Split-3D-SpGEMM does not scale well beyond 8,192 processor when squaring this matrix.

Next, we discuss strong scaling of Split-3D-SpGEMM for randomly generated square matrices whose dimensions range from  $2^{24}$  to  $2^{27}$ . Figures 5.11, 5.12a, and 5.12b show the strong scaling of multiplying two structurally similar random matrices from classes G500, ER, and SSCA, respectively. Once again, Split-3D-SpGEMM scales better when multiplying larger (e.g., scale 27) and denser matrices (G500 and ER matrices have 16 nonzeros per row, but SSCA matrices have 8 nonzeros per row). Multiplying matrices with more nonzeros per row and column is expected to yield better scalability for these matrices.

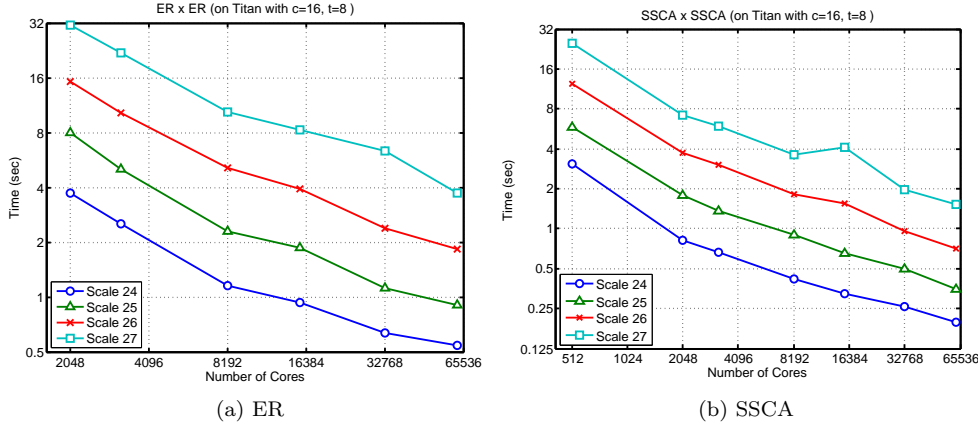


Fig. 5.12: Strong scaling of Split-3D-SpGEMM with  $c = 16, t = 8$  on Titan when multiplying two (a) ER and (b) SSCA matrices.

**5.3. Multiplication with the Restriction Operator.** Multilevel methods are widely used in the solution of numerical and combinatorial problems. Such methods construct smaller problems by successive coarsening. The simplest coarsening is graph contraction: a contraction step chooses two or more vertices in the original graph  $G$  to become a single aggregate vertex in the contracted graph  $G'$ . The edges of  $G$  that used to be incident to any of the vertices forming the aggregate become incident to the new aggregate vertex in  $G'$ . Constructing coarser representations in AMG or graph partitioning [28] is a generalized graph contraction operation. This operation can be performed by multiplying the matrix representing the original fine domain (grid, graph, or hypergraph) by the restriction operator from the left and by the transpose of the restriction from the right [25].

In our experiments, we construct the restriction matrix  $\mathbf{R}$  using distance-2 maximal independent set computation, as described by Bell et al. [7]. An independent set in a graph  $G(V, E)$  is a subset of its vertices in which no two are neighbors. A maximal independent set (MIS) is an independent set that is not a subset of any other independent set. MIS-2 is a generalization of MIS where no two vertices are distance-2 neighbors. In this scheme, each aggregate is defined by a vertex in MIS-2 and consists of the union of that vertex with its distance-1 neighbors.

The linear algebraic formulation of Luby’s randomized MIS algorithm [37] was originally described earlier [38]. Here, we generalize it to distance-2 case, which is shown in Algorithm 3 at a high level. MXV signifies matrix-vector multiplication. EWISEADD performs element-wise addition between two vectors, which amounts to a union operation among the index sets of those vectors. EWISEMULT is the element-wise multiplication, which amounts to an intersection operation among the index sets. For both EWISEADD and EWISEMULT, wherever the index sets of two vectors overlap, the values for the overlapping indices are “added” according to the binary function that is passed as the third parameter. Line 5 finds the smallest random value among a vertex’s neighbors using the semiring where scalar multiplication is overloaded with the operation that returns the second operand and the scalar addition is overloaded with the minimum operation. Line 6 extends this to find the smallest random value among the 2-hop neighborhood. Line 8 returns the new additions to MIS-2 if the

random value of the second vector (`cands`) is smaller. Line 9 removes those new additions, `newS`, from the set of candidates. The rest of the computation is self-explanatory.

---

**Algorithm 3** Pseudocode for MIS-2 computation in the language of matrices

---

**Input:**  $\mathbf{A} \in \mathbb{S}^{n \times n}$ ,  $\mathbf{cands} \in \mathbb{S}^{1 \times n}$

**Output:**  $\mathbf{mis2} \in \mathbb{S}^{1 \times n}$ : distance-2 maximal independent set, empty in the beginning

```

1: procedure MIS2( $\mathbf{A}$ ,  $\mathbf{cands}$ ,  $\mathbf{mis2}$ )
2:    $\mathbf{cands} = 1 : n$  ▷ all vertices are initially candidates
3:   while NNZ( $\mathbf{cands}$ ) > 0 do
4:     APPLY( $\mathbf{cands}$ , RAND()) ▷ generate random values
5:      $\mathbf{minadj1} \leftarrow \text{MXV}(\mathbf{A}, \mathbf{cands}, \text{SEMIRING}(\mathit{min}, \mathit{select2nd}))$ 
6:      $\mathbf{minadj2} \leftarrow \text{MXV}(\mathbf{A}, \mathbf{minadj1}, \text{SEMIRING}(\mathit{min}, \mathit{select2nd}))$ 
7:      $\mathbf{minadj} \leftarrow \text{WISEADD}(\mathbf{minadj1}, \mathbf{minadj2}, \text{MIN}())$  ▷ Union of minimums
8:      $\mathbf{newS} \leftarrow \text{WISEMULT}(\mathbf{minadj}, \mathbf{cands}, \text{IS2NDSMALLER}())$ 
9:      $\mathbf{cands} \leftarrow \text{WISEMULT}(\mathbf{cands}, \mathbf{newS}, \text{SELECT1ST}())$ 
10:     $\mathbf{newS\_adj1} \leftarrow \text{MXV}(\mathbf{A}, \mathbf{newS}, \text{SEMIRING}(\mathit{min}, \mathit{select2nd}))$ 
11:     $\mathbf{newS\_adj2} \leftarrow \text{MXV}(\mathbf{A}, \mathbf{newS\_adj1}, \text{SEMIRING}(\mathit{min}, \mathit{select2nd}))$ 
12:     $\mathbf{newS\_adj} \leftarrow \text{WISEADD}(\mathbf{newS\_adj1}, \mathbf{newS\_adj2}, \text{ANY}())$  ▷ Union of neighbors
13:     $\mathbf{cands} \leftarrow \text{WISEMULT}(\mathbf{cands}, \mathbf{newS\_adj}, \text{SELECT1ST}())$ 
14:     $\mathbf{mis2} \leftarrow \text{WISEADD}(\mathbf{mis2}, \mathbf{newS}, \text{SELECT1ST}())$  ▷ Add newS to mis2

```

---

Once the set `mis2` is computed, we construct the restriction matrix  $\mathbf{R}$  by having each column represent the union of a vertex in `mis2` with its distance-1 neighborhood. The neighborhood is calculated using another MXV operation. The remaining singletons are assigned to an aggregate randomly in order to ensure good load balance. Consequently,  $\mathbf{R}$  is of dimensions  $n \times \text{size}(\mathbf{mis2})$ .

**5.3.1. Performance of Multiplying a Matrix with the Restriction Operator.** Figure 5.13 shows the strong scaling of different variants of 2D and 3D algorithms when computing  $\mathbf{R}^T \mathbf{A}$  with `NaLuR3` matrix on Edison. Split-3D-SpGEMM with  $c = 16, t = 6$  attains  $7.5 \times$  speedup when we go from 512 cores to 32,768 cores, but other variants of 2D and 3D algorithms achieve lower speedups. Comparing the scaling of squaring `NaLuR3` from Figure 5.5, we observe moderate scalability of all variants of 2D and 3D algorithms when multiplying `NaLuR3` with the restriction matrix. This is because the number of nonzeros in  $\mathbf{R}^T \mathbf{A}$  is only 77 million, whereas  $\text{nnz}(\mathbf{A}^2) = 2.1$  billion for `NaLuR3` matrix (see Table 5.2). Hence, unlike squaring `NaLuR3`,  $\mathbf{R}^T \mathbf{A}$  computation does not have enough work to utilize thousands of cores. However, the performance gap between 2D and 3D algorithms is larger when computing  $\mathbf{R}^T \mathbf{A}$ . Figure 5.13 shows that Split-3D-SpGEMM with  $c = 16, t = 6$  runs  $8 \times$  and  $16 \times$  faster than non-threaded 2D algorithm on 512 and 32,768 cores, respectively.

Figure 5.14 shows the breakdown of runtime spent by Split-3D-SpGEMM ( $c = 16, t = 6$ ) to compute  $\mathbf{R}^T \mathbf{A}$  for (a) `n1pkkt160`, and (b) `NaLuR3` matrices on Edison. We observe that when computing  $\mathbf{R}^T \mathbf{A}$ , Split-3D-SpGEMM spends a small fraction of total runtime in the multiway merge routine. For example, on 384 cores in Figure 5.14(b), 37% of total time is spent on computation, and only about 7% of total time is spent on multiway merge. This is because  $\text{nnz}(\mathbf{R}^T \mathbf{A})$  is smaller than  $\text{nnz}(\mathbf{A})$  for `NaLuR3` matrix (also true for other matrices in Table 5.2). Therefore, in computing  $\mathbf{R}^T \mathbf{A}$ ,  $\text{nnz}(\mathbf{R}^T \mathbf{A})$  dominates the runtime of multiway merge while  $\text{nnz}(\mathbf{A})$  dominates the local multiplication, making the former less computationally intensive. Hence,

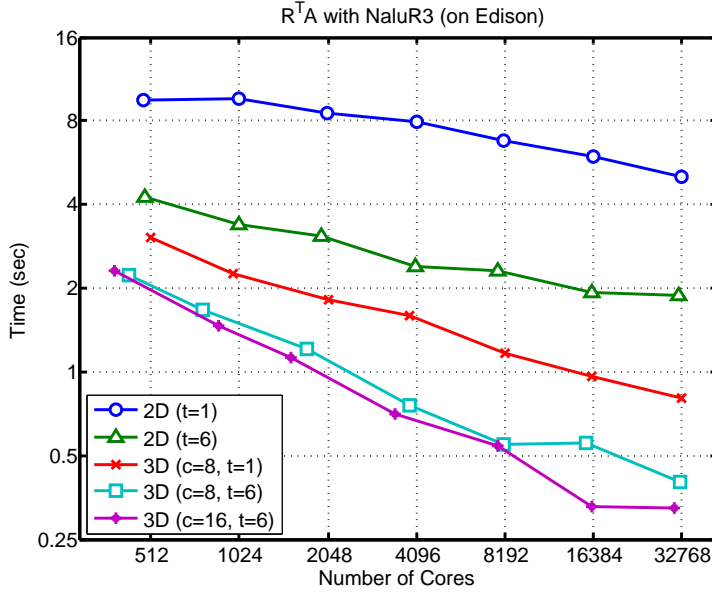


Fig. 5.13: Strong scaling of different variants of 2D and 3D algorithms to compute  $\mathbf{R}^T \mathbf{A}$  for NaluR3 matrix on Edison.

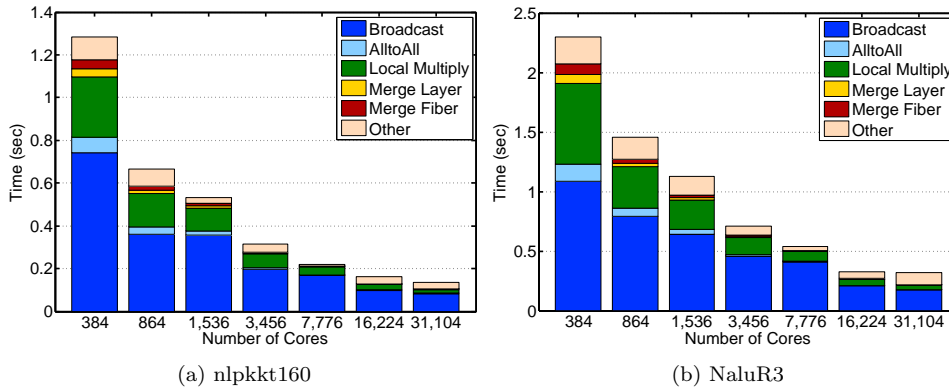


Fig. 5.14: Breakdown of runtime spent by Split-3D-SpGEMM to compute  $\mathbf{R}^T \mathbf{A}$  for (a) `nlpkkt160`, and (b) `NaluR3` matrices with  $c = 16, t = 6$  on Edison. Both communication and computation time scale well as we increase the number of cores.

despite good scaling of local multiplication, the overall runtime is dominated by communication even on lower concurrency, thereby limiting the overall scaling on tens of thousands of cores. By contrast, Split-3D-SpGEMM spends 64% of its total runtime in computation (with 40% of the total runtime spent in multiway merge) when squaring NaluR3 on 384 cores of Edison (Figure 5.14(b)). Hence, squaring of matrices shows better strong scaling than multiplying matrices with restriction operators.

Finally, Figure 5.15 shows the strong scaling of Split-3D-SpGEMM ( $c = 16, t = 6$ ) to compute  $\mathbf{R}^T \mathbf{A}$  for other real matrices. Split-3D-SpGEMM attains moderate



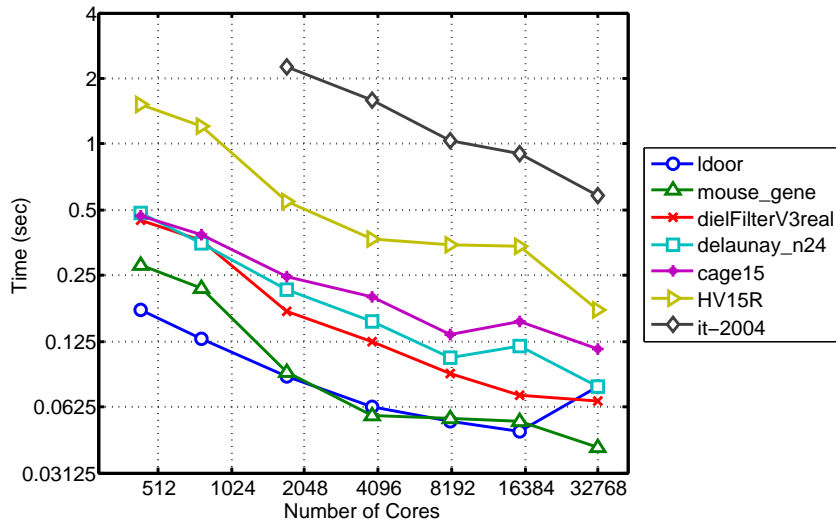


Fig. 5.15: Strong scaling of Split-3D-SpGEMM to compute  $\mathbf{R}^T \mathbf{A}$  with  $c = 16, t = 6$  for seven real matrices on Edison.

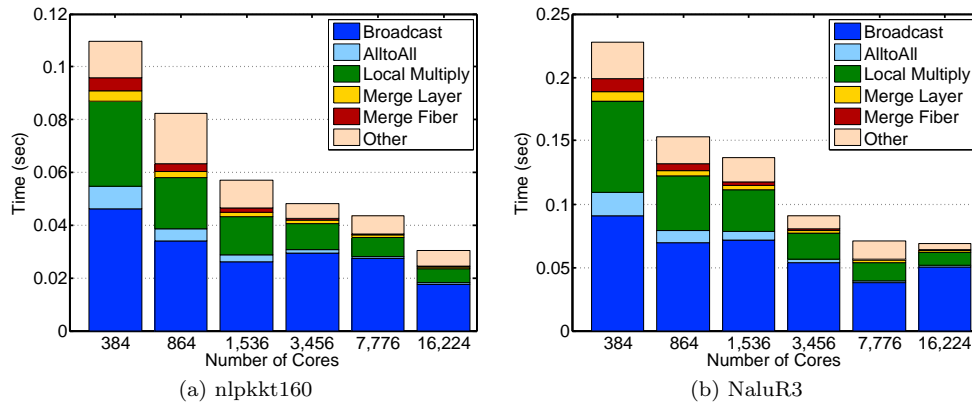


Fig. 5.16: Breakdown of runtime spent by Split-3D-SpGEMM to multiply  $\mathbf{R}^T \mathbf{A}$  and  $\mathbf{R}$  for (a) `nlpkkt160`, and (b) `NaluR3` matrices with  $c = 16, t = 6$  on Edison.

speedups of up to  $10\times$  when we got from 512 cores to 32,768 cores because of low computational intensity in the computation of  $\mathbf{R}^T \mathbf{A}$ .

**5.3.2. Performance of Multiplying  $\mathbf{R}^T \mathbf{A}$  and  $\mathbf{R}$ .** Figure 5.16 shows the scaling and breakdown of runtime spent by Split-3D-SpGEMM ( $c = 16, t = 6$ ) to multiply  $\mathbf{R}^T \mathbf{A}$  and  $\mathbf{R}$  for (a) `nlpkkt160`, and (b) `NaluR3` matrices on Edison. Even though  $(\mathbf{R}^T \mathbf{A})\mathbf{R}$  computation can still obtain limited speedups on higher concurrency, the runtimes in Figure 5.16 and the number of nonzeros in  $\mathbf{R}^T \mathbf{A}\mathbf{R}$  suggest that we might want to perform this multiplication on lower concurrency if necessary without degrading the overall performance.

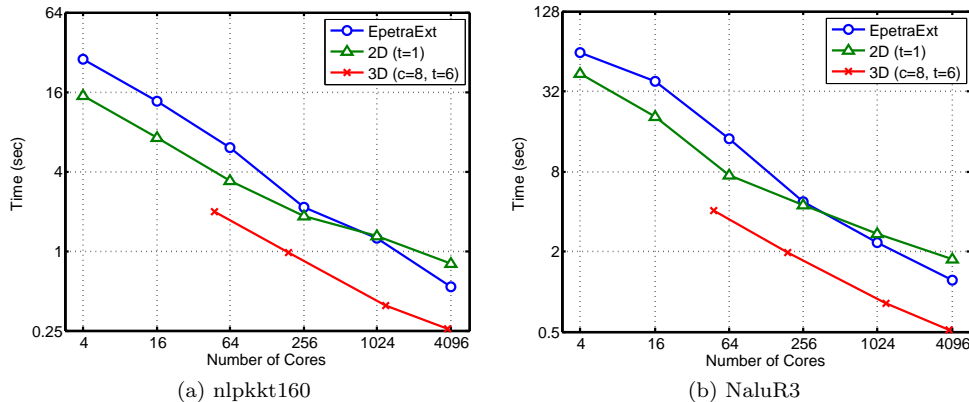


Fig. 5.17: Comparison of Trilinos’s EpetraExt package with 2D and 3D algorithms when computing  $\mathbf{AR}$  for `nlpkkt160` and `NaluR3` matrices on Edison.

**5.4. Comparison with Trilinos.** We compared the performance of our algorithms with the distributed-memory SpGEMM available in EpetraExt package of Trilinos. We observed that SpGEMM in EpetraExt runs up to  $3\times$  faster when we compute  $\mathbf{AR}$  instead of  $\mathbf{R}^T\mathbf{A}$ , especially on lower concurrency. Hence, we only consider the runtime of  $\mathbf{AR}$  so that we compare against the best configuration of EpetraExt. By contrast, our 2D and 3D algorithms are less sensitive to the order of matrix multiplication with less than  $1.5\times$  performance improvement in computing  $\mathbf{AR}$  over  $\mathbf{R}^T\mathbf{A}$ . We use a random partitioning of rows to processors for EpetraExt runs.

Figure 5.17 shows the strong scaling of EpetraExt’s SpGEMM implementation and our 2D/3D algorithms when computing  $\mathbf{AR}$  on Edison. On low concurrency, EpetraExt runs slower than the 2D algorithm, but the former eventually outperforms the latter on higher concurrency. However, on all concurrencies, the 3D algorithm with  $c = 8, t = 6$  runs at least twice as fast as EpetraExt for these two matrices. We note that these matrices are structured with good separators where 1D decomposition used in EpetraExt usually performs better. However, given the limitations of 1D decomposition for matrices without good separators, EpetraExt is not expected to perform well for graphs with power-law distributions [9]. We have tried scale 24 Graph500 matrices in EpetraExt, but received segmentation fault in I/O. We also tried other larger matrices, but EpetraExt could not finish reading the matrices from files in 24 hours, the maximum allocation limit for small jobs in Edison. Hence, we compare with EpetraExt on problems that it excels (AMG style reduction with matrices having good separators) and even there our 3D algorithm does comparably better. We could have separated the diagonal for better scaling performance [14], but we decided not to as it would break the “black box” nature of our algorithm.

**6. Conclusions and Future Work.** We presented the first implementation of the 3D parallel formulation of sparse matrix-matrix multiplication (SpGEMM). Our implementation exploits inter-node parallelism within a third processor grid dimension as well as thread-level parallelism within the node. It achieves higher performance compared to other available formulations of distributed-memory SpGEMM, without compromising flexibility in the numbers of processors that can be utilized. In particular, by varying the third processor dimension as well as the number of threads, one

can run our algorithm on many processor counts.

The percentage of time spent in communication (data movement) is significantly lower in our new implementation compared to a 2D implementation. This is advantageous for multiple reasons. First, the bandwidth for data movement is expected to increase at a slower rate than other system components, providing a future bottleneck. Second, communication costs more energy than computation [32, Figure 5]. Lastly, communication can be hidden by overlapping it with local computation, up to the time it takes to do the local computation. For example, up to 100% performance increase can be realized with overlapping if the communication costs 50% of overall time. However, if the communication costs 80% of the time, then overlapping can only increase performance by up to 25%. Overlapping communication with computation as well as exploiting task-based programming models are subject to future work.

Our 3D implementation inherits many desirable properties of the 2D matrix decomposition, such as resiliency against matrices with skewed degree distribution that are known to be very challenging for traditional 1D distributions and algorithms. However, the 3D formulation also avoids some of the pitfalls of 2D algorithms, such as their relatively poor performance on structured matrices (due to load imbalance that occurs on the processor on the diagonal), by exploiting parallelism along the third dimension. This enabled our algorithm to beat a highly-tuned 1D implementation (the new EpetraExt) on structured matrices, without resorting to techniques such as matrix splitting that were previously required of the 2D algorithm for mitigating the aforementioned load imbalance [14].

Our experimental results indicate that at large concurrencies, performance of the inter-node communication collectives becomes the determinant factor in overall performance. Even though work on the scaling of collectives on subcommunicators is under way, we believe that the effect of simultaneous communication on several subcommunicators are not well studied and should be the subject of further research.

**Acknowledgments.** We sincerely thank the anonymous reviewers whose feedback greatly improved the presentation and clarity of this paper.

This material is based upon work supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, Applied Mathematics program under contract number DE-AC02-05CH11231.

This research was supported in part by an appointment to the Sandia National Laboratories Truman Fellowship in National Security Science and Engineering, sponsored by Sandia Corporation (a wholly owned subsidiary of Lockheed Martin Corporation) as Operator of Sandia National Laboratories under its U.S. Department of Energy Contract No. DE-AC04-94AL85000.

The research of some of the authors was supported by the U.S. Department of Energy Office of Science, Office of Advanced Scientific Computing Research, Applied Mathematics program under award de-sc0010200, by the U.S. Department of Energy Office of Science, Office of Advanced Scientific Computing Research, X-Stack program under awards de-sc0008699, de-sc0008700, and AC02-05CH11231, and by DARPA award HR0011-12-2-0016, with contributions from Intel, Oracle, and MathWorks.

Research is supported by grants 1878/14, and 1901/14 from the Israel Science Foundation (founded by the Israel Academy of Sciences and Humanities) and grant 3-10891 from the Ministry of Science and Technology, Israel. Research is also supported by the Einstein Foundation and the Minerva Foundation. This work was supported by the HUJI Cyber Security Research Center in conjunction with the Israel National Cyber Bureau in the Prime Minister's Office. This paper is supported

by the Intel Collaborative Research Institute for Computational Intelligence (ICRI-CI). This research was supported by a grant from the United States-Israel Binational Science Foundation (BSF), Jerusalem, Israel.

This research used resources of the National Energy Research Scientific Computing Center, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231, and resources of the Oak Ridge Leadership Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725.

## REFERENCES

- [1] Graph500 benchmark. [www.graph500.org](http://www.graph500.org).
- [2] SSCA benchmark. <http://www.graphanalysis.org/benchmark/>.
- [3] Kadir Akbudak and Cevdet Aykanat. Simultaneous input and output matrix partitioning for outer-product-parallel sparse matrix-matrix multiplication. *SIAM Journal on Scientific Computing*, 36(5):C568–C590, 2014.
- [4] Ariful Azad, Aydın Buluç, and John R Gilbert. Parallel triangle counting and enumeration using matrix algebra. In *Proceedings of the IPDPSW, Workshop on Graph Algorithm Building Blocks (GABB)*, 2015.
- [5] Grey Ballard, Aydın Buluç, James Demmel, Laura Grigori, Benjamin Lipshitz, Oded Schwartz, and Sivan Toledo. Communication optimal parallel multiplication of sparse random matrices. In *SPAA 2013: The 25th ACM Symposium on Parallelism in Algorithms and Architectures*, Montreal, Canada, 2013.
- [6] Grey Ballard, Alex Druinsky, Nicholas Knight, and Oded Schwartz. Brief announcement: Hypergraph partitioning for parallel sparse matrix-matrix multiplication. In *Proceedings of the 27th ACM on Symposium on Parallelism in Algorithms and Architectures (SPAA)*, pages 86–88, 2015.
- [7] Nathan Bell, Steven Dalton, and Luke N Olson. Exposing fine-grained parallelism in algebraic multigrid methods. *SIAM Journal on Scientific Computing*, 34(4):C123–C152, 2012.
- [8] Nicolas Bock and Matt Challacombe. An optimized sparse approximate matrix multiply for matrices with decay. *SIAM Journal on Scientific Computing*, 35(1):C72–C98, 2013.
- [9] Erik G Boman, Karen D Devine, and Sivasankaran Rajamanickam. Scalable matrix computations on large scale-free graphs using 2d graph partitioning. In *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*, page 50. ACM, 2013.
- [10] Urban Borštnik, Joost VandeVondele, Valéry Weber, and Jürg Hutter. Sparse matrix multiplication: The distributed block-compressed sparse row library. *Parallel Computing*, 40(5):47–58, 2014.
- [11] Aydın Buluç and John R. Gilbert. Challenges and advances in parallel sparse matrix-matrix multiplication. In *ICPP'08: Proc. of the Intl. Conf. on Parallel Processing*, pages 503–510, Portland, Oregon, USA, 2008. IEEE Computer Society.
- [12] Aydın Buluç and John R. Gilbert. On the representation and multiplication of hypersparse matrices. In *IPDPS'08: Proceedings of the IEEE International Symposium on Parallel & Distributed Processing*, pages 1–11. IEEE Computer Society, 2008.
- [13] Aydın Buluç and John R. Gilbert. The Combinatorial BLAS: Design, implementation, and applications. *International Journal of High Performance Computing Applications (IJHPCA)*, 25(4):496–509, 2011.
- [14] Aydın Buluç and John R. Gilbert. Parallel sparse matrix-matrix multiplication and indexing: Implementation and experiments. *SIAM Journal of Scientific Computing (SISC)*, 34(4):170 – 191, 2012.
- [15] Aydın Buluç, John R. Gilbert, and Viral B. Shah. Implementing Sparse Matrices for Graph Algorithms. In Jeremy Kepner and John R. Gilbert, editors, *Graph Algorithms in the Language of Linear Algebra*. SIAM, 2011.
- [16] Deepayan Chakrabarti, Yiping Zhan, and Christos Faloutsos. R-MAT: A recursive model for graph mining. In Michael W. Berry, Umeshwar Dayal, Chandrika Kamath, and David B. Skillicorn, editors, *SDM*. SIAM, 2004.
- [17] Ernie Chan, Marcel Heimlich, Avi Purkayastha, and Robert A. van de Geijn. Collective communication: theory, practice, and experience. *Concurrency and Computation: Practice*

- and Experience*, 19(13):1749–1783, 2007.
- [18] Steven Dalton, Luke Olsen, and Nathan Bell. Optimizing sparse matrix-matrix multiplication for the GPU. *ACM Transactions on Mathematical Software*, 41(4), to appear.
  - [19] Timothy A. Davis. *Direct Methods for Sparse Linear Systems*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2006.
  - [20] Timothy A Davis and Yifan Hu. The university of florida sparse matrix collection. *ACM Transactions on Mathematical Software (TOMS)*, 38(1):1, 2011.
  - [21] E. Dekel, D. Nassimi, and S. Sahni. Parallel matrix and graph algorithms. *SIAM Journal on Computing*, 10(4):657–675, 1981.
  - [22] Edison website. <http://www.nersc.gov/users/computational-systems/edison>.
  - [23] R. A. Van De Geijn and J. Watts. SUMMA: Scalable universal matrix multiplication algorithm. *Concurrency: Practice and Experience*, 9(4):255–274, 1997.
  - [24] John R. Gilbert, Cleve Moler, and Robert Schreiber. Sparse matrices in Matlab: Design and implementation. *SIAM Journal of Matrix Analysis and Applications*, 13(1):333–356, 1992.
  - [25] John R. Gilbert, Steve Reinhardt, and Viral B. Shah. A unified framework for numerical and combinatorial computing. *Computing in Science and Engineering*, 10(2):20–25, 2008.
  - [26] Felix Gremse, Andreas Hoffer, Lars Ole Schwen, Fabian Kiessling, and Uwe Naumann. Gpu-accelerated sparse matrix-matrix multiplication by iterative row merging. *SIAM Journal on Scientific Computing*, 37(1):C54–C71, 2015.
  - [27] Fred G. Gustavson. Two fast algorithms for sparse matrices: Multiplication and permuted transposition. *ACM Transactions on Mathematical Software*, 4(3):250–269, 1978.
  - [28] Bruce Hendrickson and Robert Leland. A multilevel algorithm for partitioning graphs. In *Supercomputing '95: Proceedings of the 1995 ACM/IEEE conference on Supercomputing*, page 28, New York, NY, USA, 1995. ACM.
  - [29] Michael A. Heroux, Roscoe A. Bartlett, Vicki E. Howle, Robert J. Hoekstra, Jonathan J. Hu, Tamara G. Kolda, Richard B. Lehoucq, Kevin R. Long, Roger P. Pawlowski, Eric T. Phipps, Andrew G. Salinger, Heidi K. Thornquist, Ray S. Tuminaro, James M. Willenbring, Alan Williams, and Kendall S. Stanley. An overview of the Trilinos project. *ACM Trans. Math. Softw.*, 31(3):397–423, 2005.
  - [30] D. Irony, S. Toledo, and A. Tiskin. Communication lower bounds for distributed-memory matrix multiplication. *J. Parallel Distrib. Comput.*, 64(9):1017–1026, 2004.
  - [31] J. Kepner and J.R. Gilbert, editors. *Graph Algorithms in the Language of Linear Algebra*. SIAM, Philadelphia, 2011.
  - [32] Peter Kogge and John Shalf. Exascale computing trends: Adjusting to the. *Computing in Science & Engineering*, 15(6):16–26, 2013.
  - [33] Walter Kohn. Density functional and density matrix method scaling linearly with the number of atoms. *Physical Review Letters*, 76(17):3168, 1996.
  - [34] Vijay P. Kumar and Anshul Gupta. Analyzing scalability of parallel algorithms and architectures. *Journal of parallel and distributed computing*, 22(3):379–391, 1994.
  - [35] Paul Lin, Matthew Bettencourt, Stefan Domino, Travis Fisher, Mark Hoemmen, Jonathan Hu, Eric Phipps, Andrey Prokopenko, Sivasankaran Rajamanickam, Christopher Siefert, et al. Towards extreme-scale simulations for low mach fluids with second-generation trilinos. *Parallel Processing Letters*, 24(04):1442005, 2014.
  - [36] Weifeng Liu and Brian Vinter. A framework for general sparse matrix-matrix multiplication on gpus and heterogeneous processors. *Journal of Parallel and Distributed Computing*, 2015. (to appear).
  - [37] Michael Luby. A simple parallel algorithm for the maximal independent set problem. *SIAM Journal on Computing*, 15(4):1036, 1986.
  - [38] Adam Lugowski, Shoaib Kamil, Aydın Buluç, Samuel Williams, Erika Duriakova, Leonid Oliker, Armando Fox, and John R Gilbert. Parallel processing of filtered queries in attributed semantic graphs. *Journal of Parallel and Distributed Computing*, 79:115–131, 2015.
  - [39] R. van de Geijn M. Schatz, J. Poulson. Scalable universal matrix multiplication algorithms: 2D and 3D variations on a theme. Technical report, UT Austin, 2013.
  - [40] Tim Mattson, David Bader, Jonathan Berry, Aydın Buluç, Jack Dongarra, Christos Faloutsos, John Feo, Jeremy Gilbert, Joseph Gonzalez, Bruce Hendrickson, et al. Standards for graph algorithm primitives. In *High Performance Extreme Computing Conference (HPEC)*, pages 1–2. IEEE, 2013.
  - [41] Michael McCourt, Barry Smith, and Hong Zhang. Sparse matrix-matrix products executed through coloring. *SIAM Journal on Matrix Analysis and Applications*, 36(1):90–109, 2015.
  - [42] Md Mostofa Ali Patwary, Nadathur Rajagopalan Satish, Narayanan Sundaram, Jongsoo Park, Michael J Anderson, Satya Gautam Vadlamudi, Dipankar Das, Sergey G Pudov, Vadim O Pirogov, and Pradeep Dubey. Parallel efficient sparse matrix-matrix multiplication on

- multicore platforms. In *High Performance Computing*, pages 48–57. Springer, 2015.
- [43] Viral B. Shah. *An Interactive System for Combinatorial Scientific Computing with an Emphasis on Programmer Productivity*. PhD thesis, University of California, Santa Barbara, June 2007.
- [44] Edgar Solomonik and James Demmel. Communication-optimal parallel 2.5 d matrix multiplication and lu factorization algorithms. In *Euro-Par 2011 Parallel Processing*, pages 90–109. Springer, 2011.
- [45] Titan website. <https://www.olcf.ornl.gov/titan/>.
- [46] Leslie G Valiant. Optimally universal parallel computers. In *Opportunities and Constraints of Parallel Computing*, pages 155–158. Springer, 1989.
- [47] Stijn van Dongen. *Graph Clustering by Flow Simulation*. PhD thesis, University of Utrecht, 2000.
- [48] Ichitaro Yamazaki and Xiaoye Li. On techniques to improve robustness and scalability of a parallel hybrid linear solver. In *High Performance Computing for Computational Science VECPAR*, pages 421–434, 2010.
- [49] Raphael Yuster and Uri Zwick. Detecting short directed cycles using rectangular matrix multiplication and dynamic programming. In *SODA '04: Proceedings of the fifteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 254–260, 2004.